

GPT et al.

Generating Texts with Transformer-Based Large Language Models



CEEHACKS



matfyz



Charles University
Faculty of Mathematics and Physics
Institute of Formal and Applied Linguistics



unless otherwise stated

Motivation: Transformer-based models

- Vaswani+, 2017: Attention is All you Need (NIPS)
 - State-of-the-art for most Natural Language Processing (NLP) tasks
 - Machine Translation human parity (*)
 - Passing the Turing Test (*)
- OpenAI: GPT-1 (2018) ... GPT-4 (2023)
 - Large Language Models (LLM)
 - Generating texts, program code, processing data
 - Integration into tools (Copilot, Bing, Duolingo...)
 - Passing various highschool & university exams
 - Generating theatre play scripts (THEaiTRE)
- Multimodal models
 - Vision Transformers
 - GPT-4: input = text + image
 - PaLM-E: multimodal + embodiment

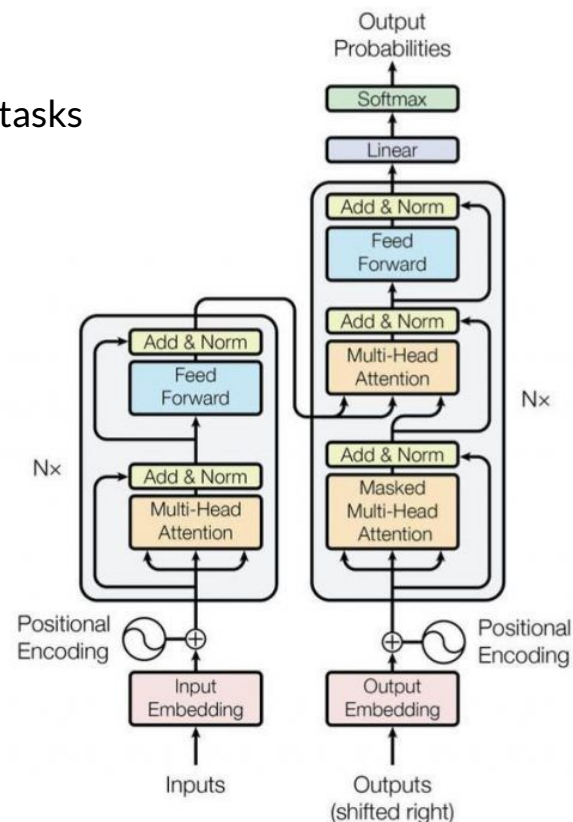


Figure 1: The Transformer - model architecture.



Let's watch a sample of "AI: When a Robot Writes a Play"!

DIRECTOR DANIEL HRBEK

**AI: WHEN A ROBOT
WRITES A PLAY**



Outline

- Motivation
- Ground for Transformers
- Transformer architecture
- Models and applications

Running example: Machine Translation

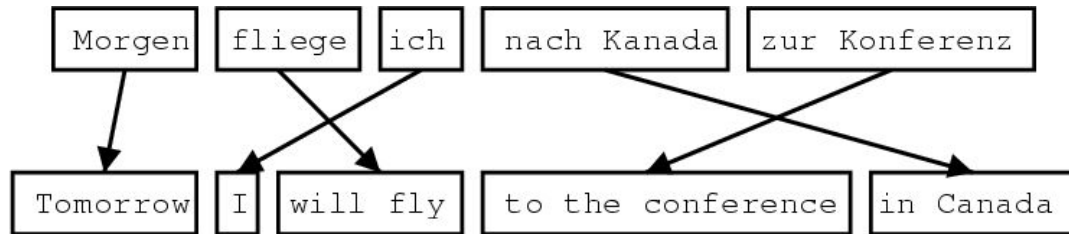
- All human beings are born free and equal in dignity and rights.



- Všichni lidé rodí se svobodní a sobě rovní co do důstojnosti a práv.

Pre-neural Machine Translation

- Human translation
 - All human beings are born free and equal in dignity and rights.
 - Všichni lidé rodí se svobodní a sobě rovní co do důstojnosti a práv.
- Word-based
 - All human beings are born free and equal in dignity and rights .
 - Všechna lidská stvoření jsou zrozena svobodná a rovná v důstojnosti a právech .
- Phrase-based
 - All human beings are born free and equal in dignity and rights .
 - Všichni lidé rodí se svobodní a sobě rovní v důstojnosti a právech .
- Components
 - translation model = simple frequency tables → phrase translation suggestions
 - reordering model (word alignment)
 - language model → combines phrase translations into a fluent sentence



Yet another example: Language Modelling

- Discriminative task: Which sentence is better (is more probable)?
 - All human beings are born free and equal in dignity and rights.
 - All human beings are born unmarried and equal in dignity and rights.
 - All humans beings are born free and equal in dignity and rights.
 - Free and equal in dignity and rights all human beings are born.

Yet another example: Language Modelling

- Discriminative task: Which sentence is better (is more probable)?
 - All human beings are born free and equal in dignity and rights.
 - All human beings are born unmarried and equal in dignity and rights.
 - All humans beings are born free and equal in dignity and rights.
 - Free and equal in dignity and rights all human beings are born.
- Generative task: Which word should follow (is more probable)?
 - I woke up in the morning and went to the...

Yet another example: Language Modelling

- Discriminative task: Which sentence is better (is more probable)?
 - All human beings are born free and equal in dignity and rights.
 - All human beings are born unmarried and equal in dignity and rights.
 - All humans beings are born free and equal in dignity and rights.
 - Free and equal in dignity and rights all human beings are born.
- Generative task: Which word should follow (is more probable)?
 - I woke up in the morning and went to the...
 - kitchen
 - bathroom
 - cinema
 - horse

Yet another example: Language Modelling

- Discriminative task: Which sentence is better (is more probable)?
 - All human beings are born free and equal in dignity and rights.
 - All human beings are born unmarried and equal in dignity and rights.
 - All humans beings are born free and equal in dignity and rights.
 - Free and equal in dignity and rights all human beings are born.
- Generative task: Which word should follow (is more probable)?
 - I woke up in the morning and went to the...
 - kitchen
 - bathroom
 - cinema
 - horse
- Semantic tasks:
 - What is the meaning of the sentence?
 - How similar are meanings of two given sentences?
 - ...

Pre-neural Language Models

- Discrimination/generation
 - I woke up in the morning and went to the... kitchen/bathroom/cinema/horse
- N-gram language models (e.g. 4-grams)
 - How often are words ABC followed by word **D** (in first 100M lines of Wikipedia)?
 - “went to the bathroom” > “went to the horse”?
 - “went to the cinema” > “went to the kitchen”?
 - simple frequency tables, hard conditioning, limited context, no concept of word similarity
 - not usable for semantic tasks

Pre-neural Language Models

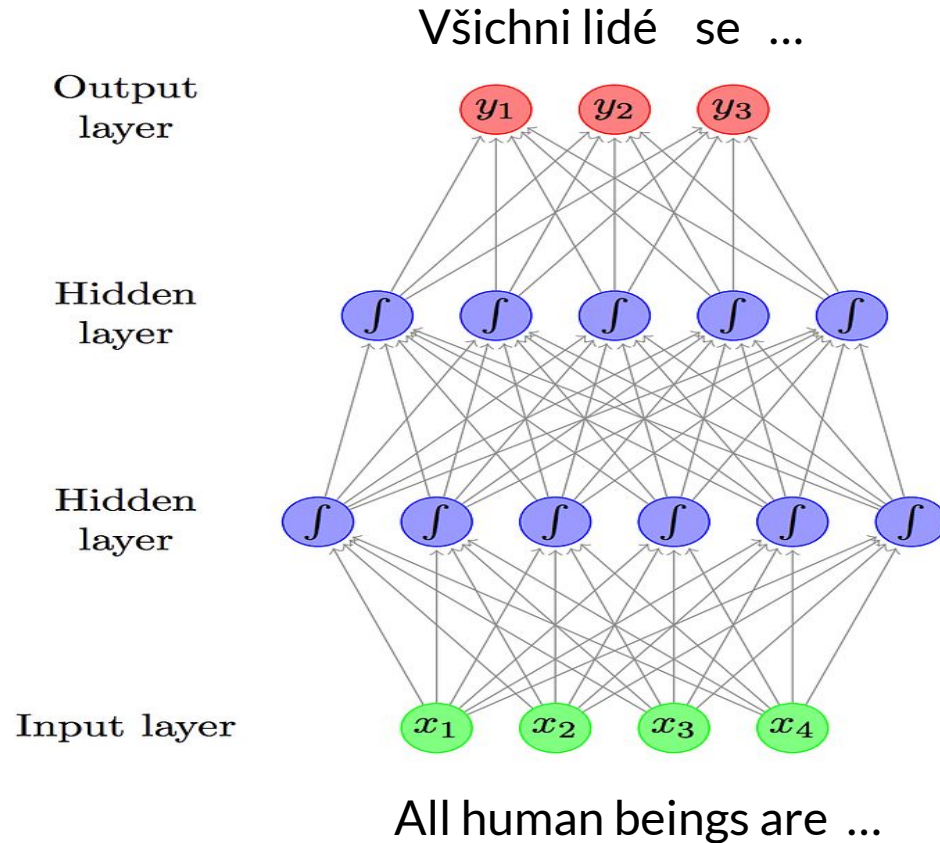
- Discrimination/generation
 - I woke up in the morning and went to the... kitchen/bathroom/cinema/horse
- N-gram language models (e.g. 4-grams)
 - How often are words ABC followed by word **D** (in first 100M lines of Wikipedia)?
 - “went to the bathroom” > “went to the horse”? 42 > 4
 - “went to the cinema” > “went to the kitchen”?
 - simple frequency tables, hard conditioning, limited context, no concept of word similarity
 - not usable for semantic tasks

Pre-neural Language Models

- Discrimination/generation
 - I woke up in the morning and went to the... kitchen/bathroom/cinema/horse
- N-gram language models (e.g. 4-grams)
 - How often are words ABC followed by word **D** (in first 100M lines of Wikipedia)?
 - “went to the bathroom” > “went to the horse”? 42 > 4
 - “went to the cinema” > “went to the kitchen”? 25 > 14
 - simple frequency tables, hard conditioning, limited context, no concept of word similarity
 - not usable for semantic tasks

Ground for Transformers

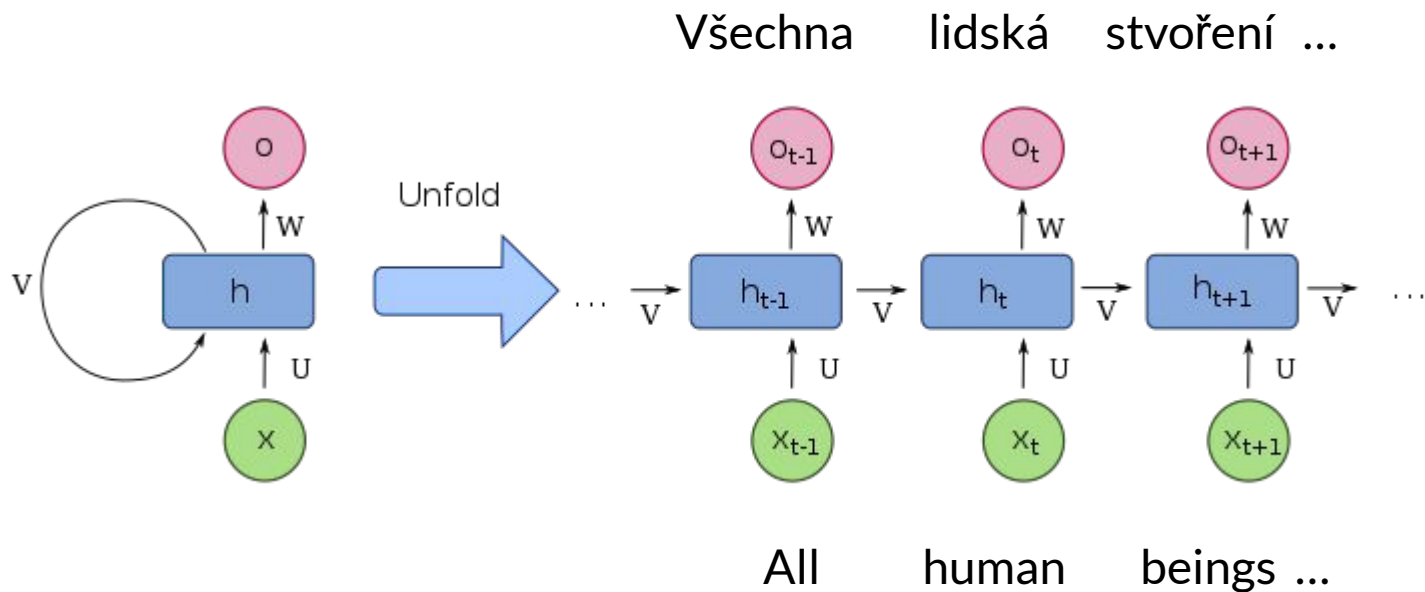
Feed Forward Neural Network



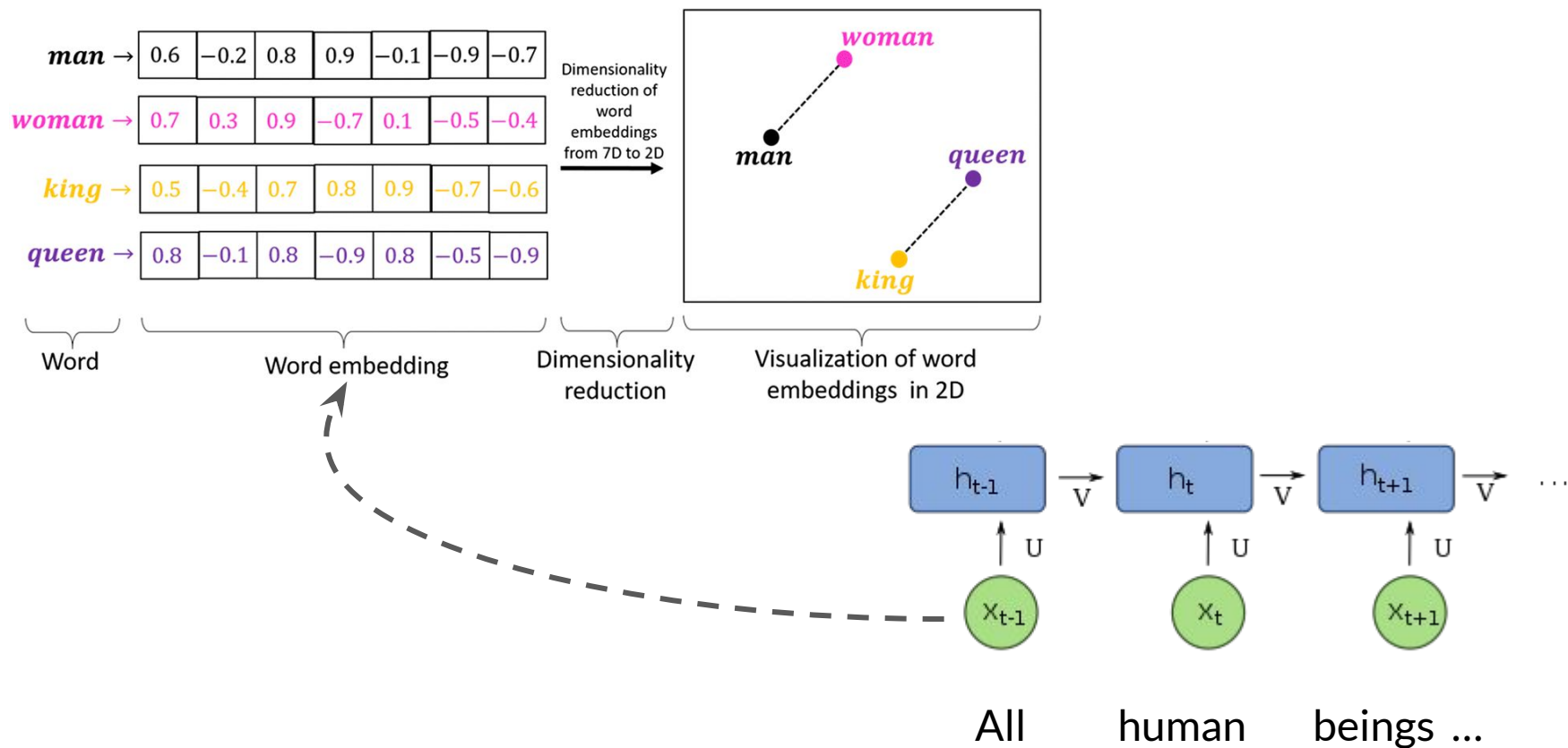
Gianp,

<https://medium.com/mllearning-ai/feedforward-neural-networks-multi-layers-preceptors-mlps-1bea7ff11e07>

Recurrent Neural Network (RNN)



Word encoding: word embeddings (e.g. word2vec)

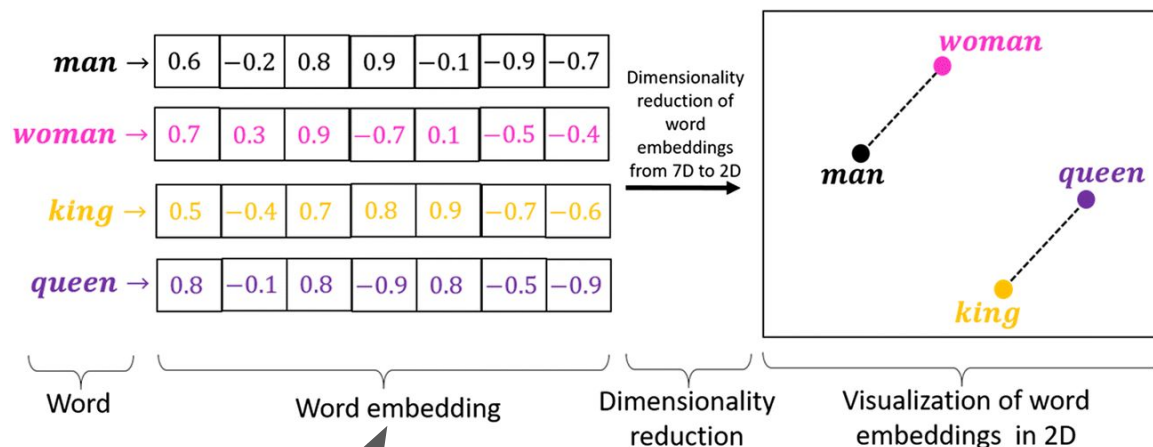


Siddharth M,

<https://www.analyticsvidhya.com/blog/2021/07/>

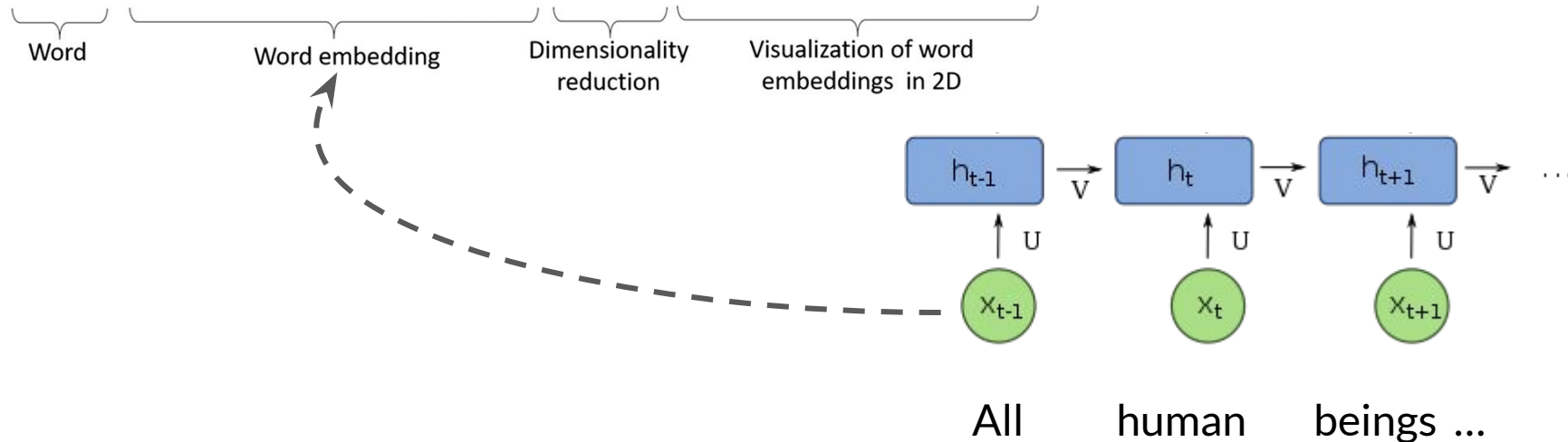
[feature-extraction-and-embeddings-in-nlp-a-beginners-guide-to-understand-natural-language-processing/](https://www.analyticsvidhya.com/blog/2021/07/feature-extraction-and-embeddings-in-nlp-a-beginners-guide-to-understand-natural-language-processing/)

Word encoding: word embeddings (e.g. word2vec)



...and actually typically also split to subwords:

platypus → *plat- ypu- s*

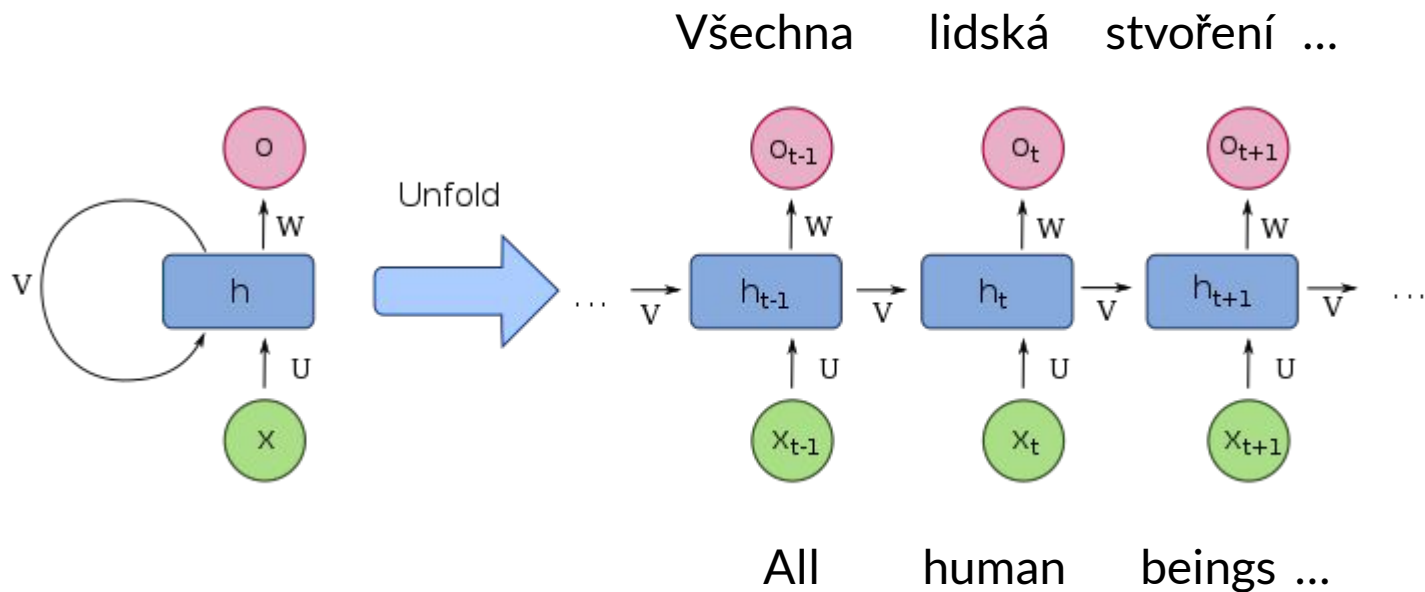


Siddharth M,

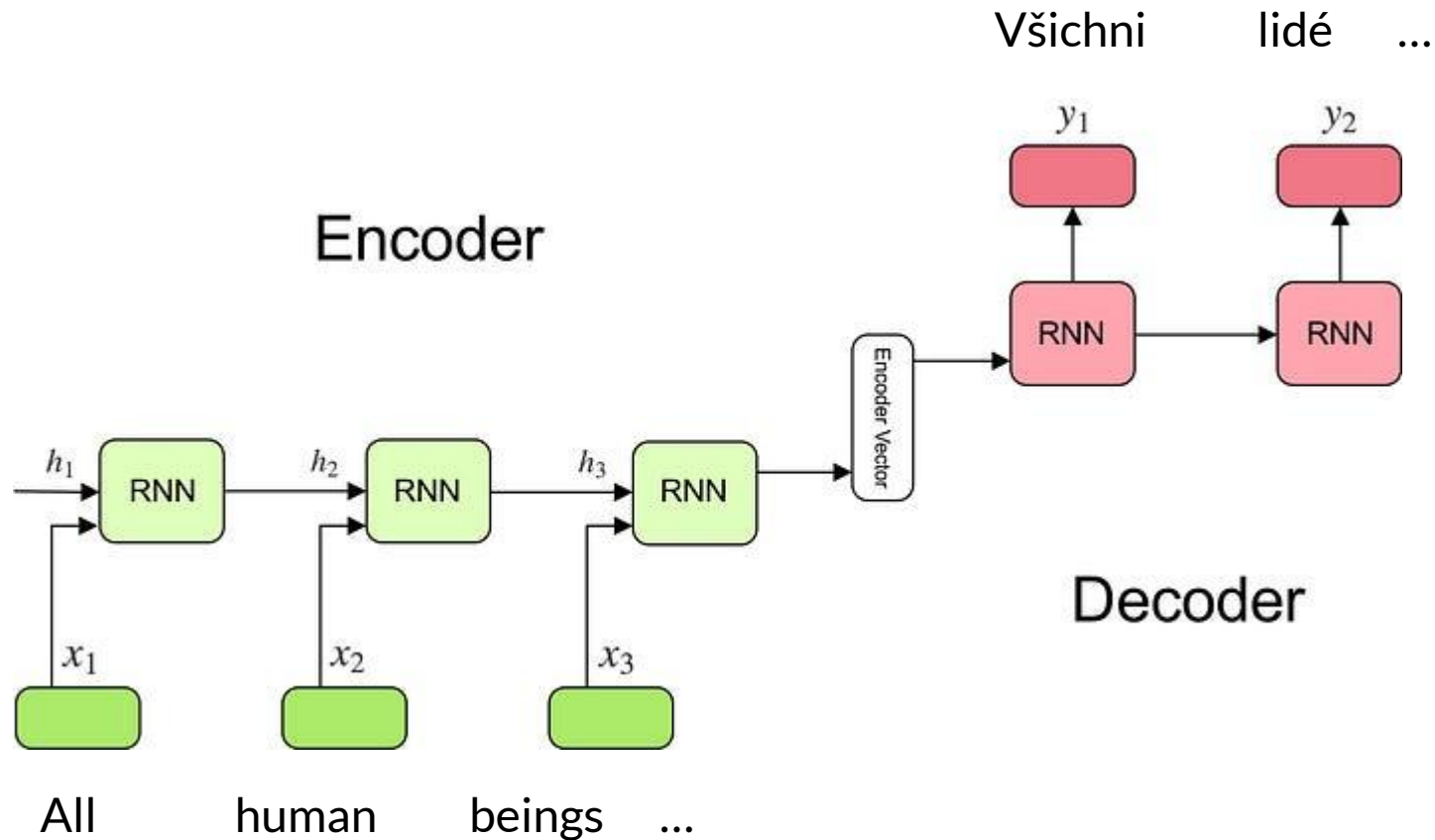
<https://www.analyticsvidhya.com/blog/2021/07/>

[feature-extraction-and-embeddings-in-nlp-a-beginners-guide-to-understand-natural-language-processing/](https://www.analyticsvidhya.com/blog/2021/07/feature-extraction-and-embeddings-in-nlp-a-beginners-guide-to-understand-natural-language-processing/)

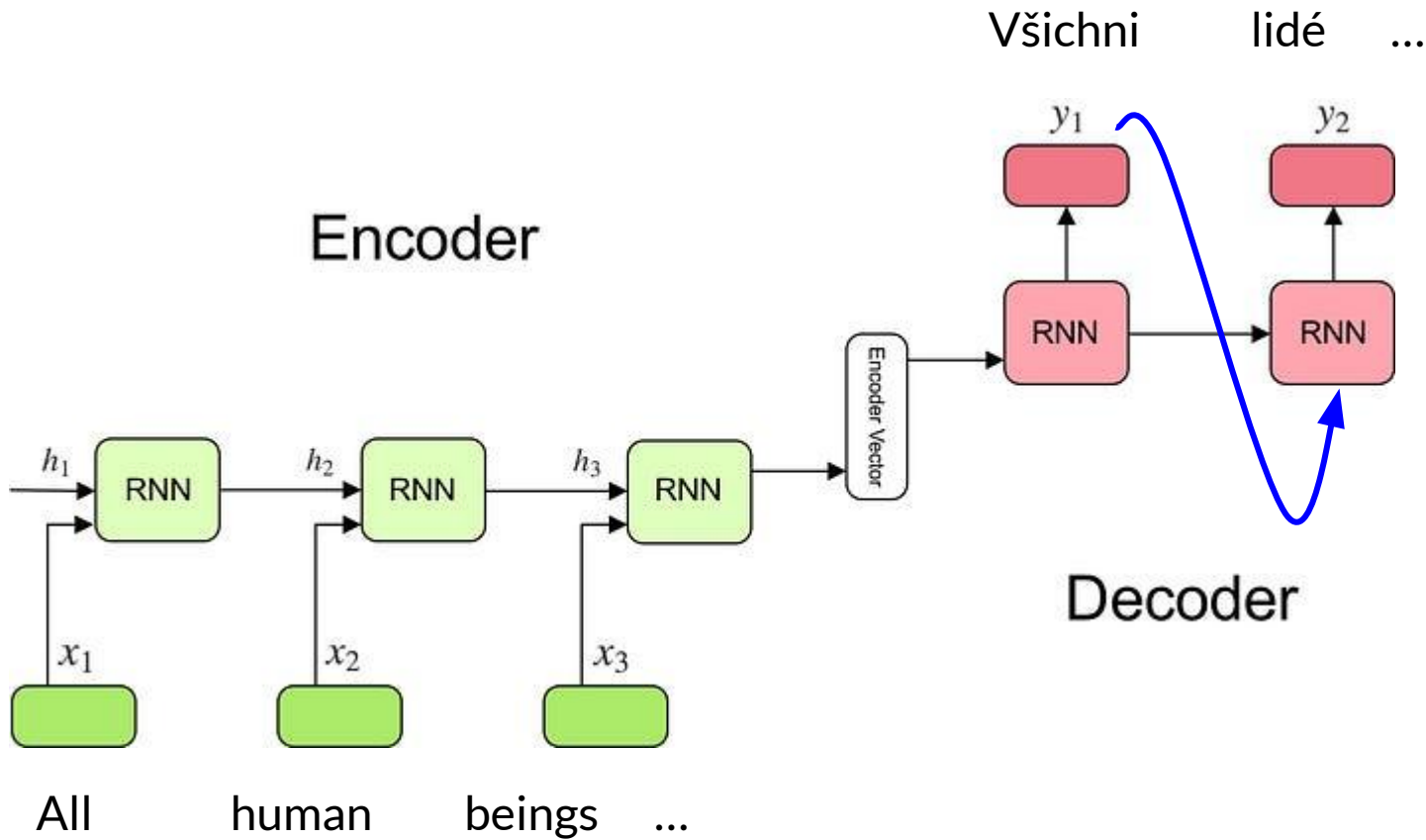
Recurrent Neural Network (RNN)



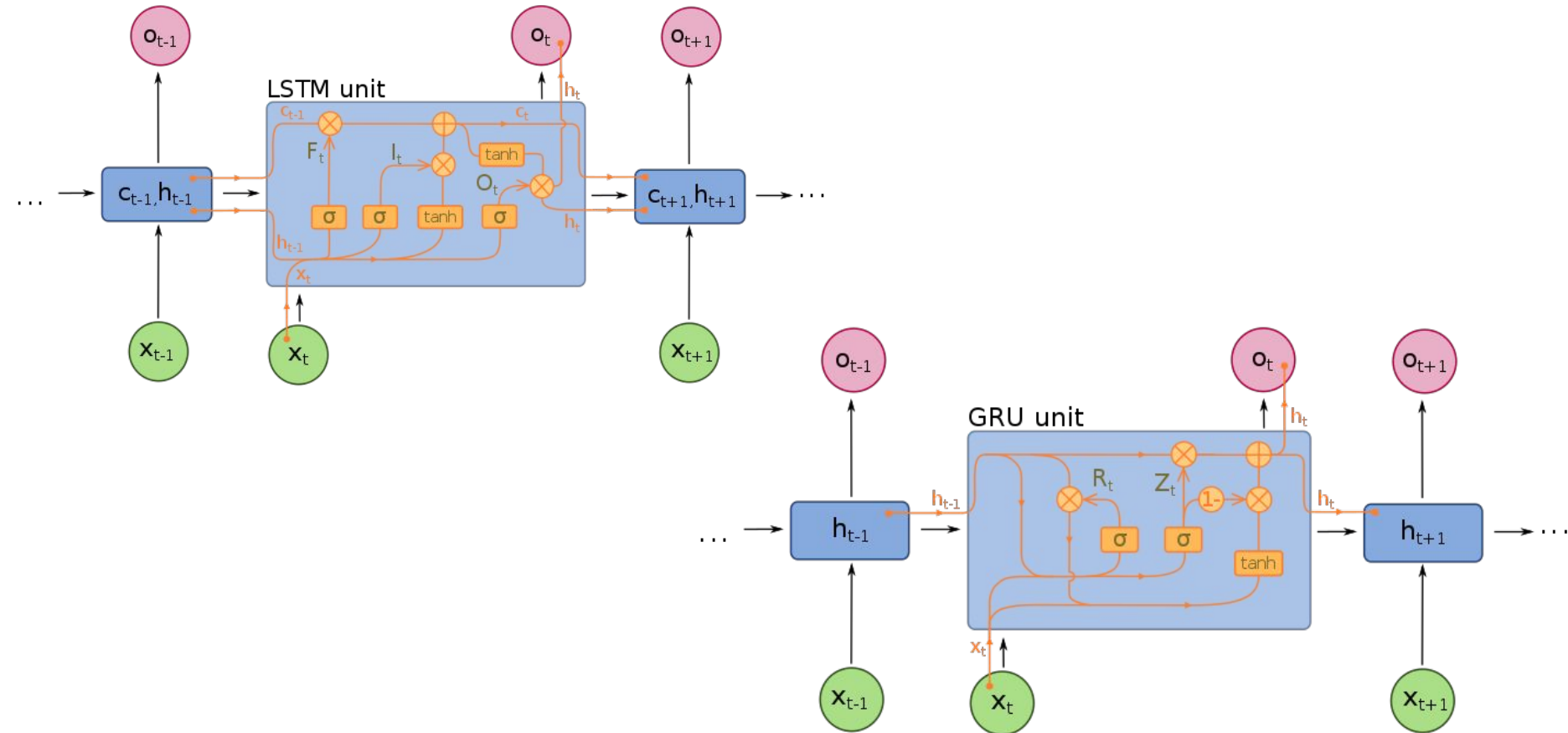
Encoder-decoder (sequence-to-sequence)



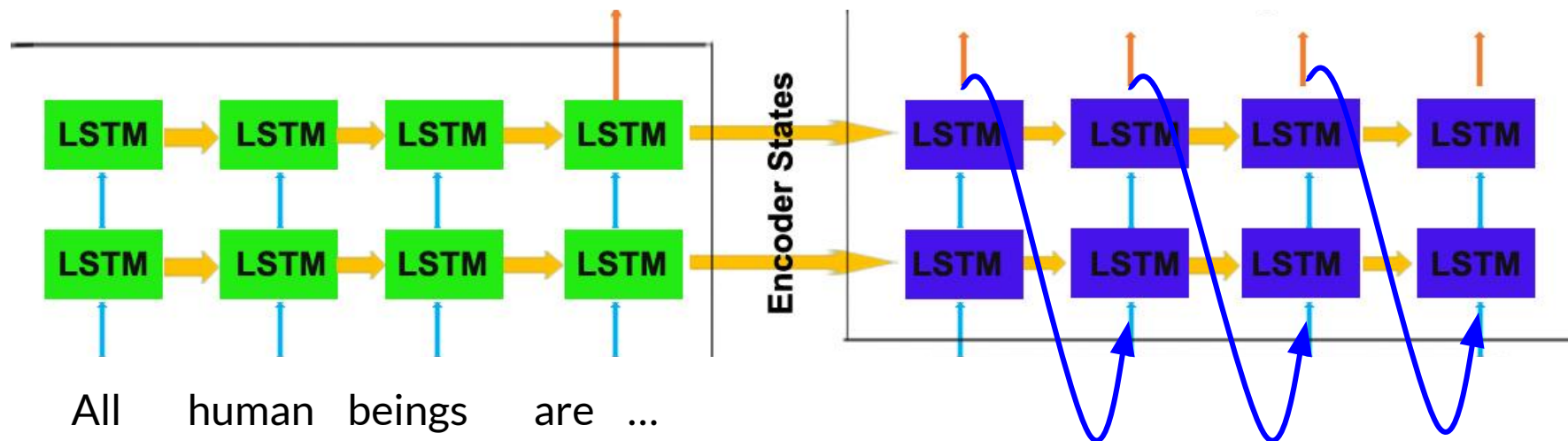
Based on previous output → Autoregressive Encoder-decoder



Vanishing/Exploding Gradient → RNNs with memory: LSTM, GRU



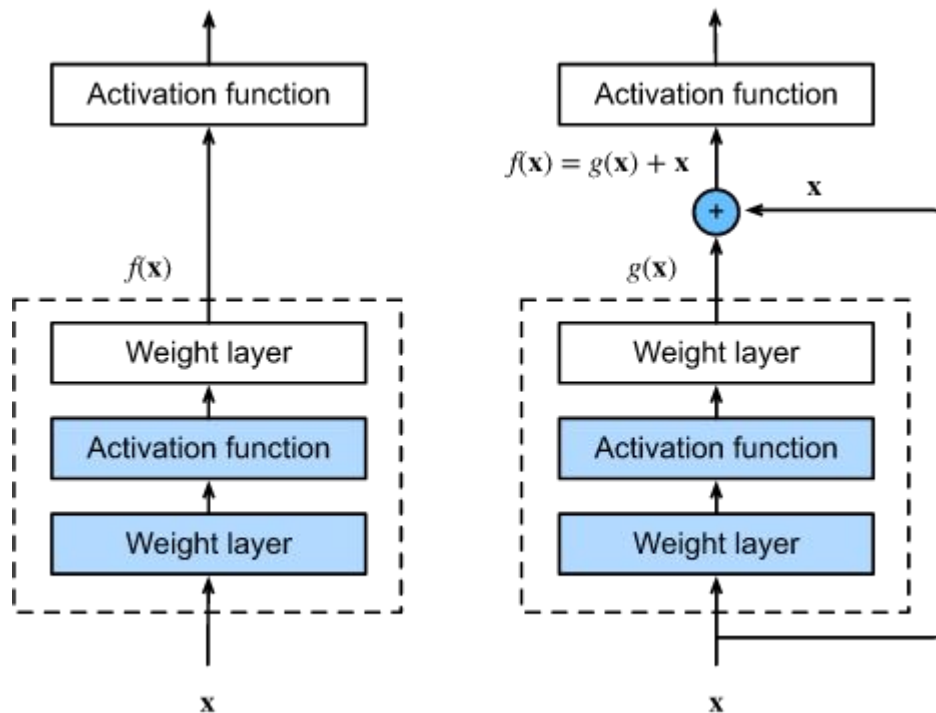
Encode larger context → Stacked RNN



S. Jagadeesh,

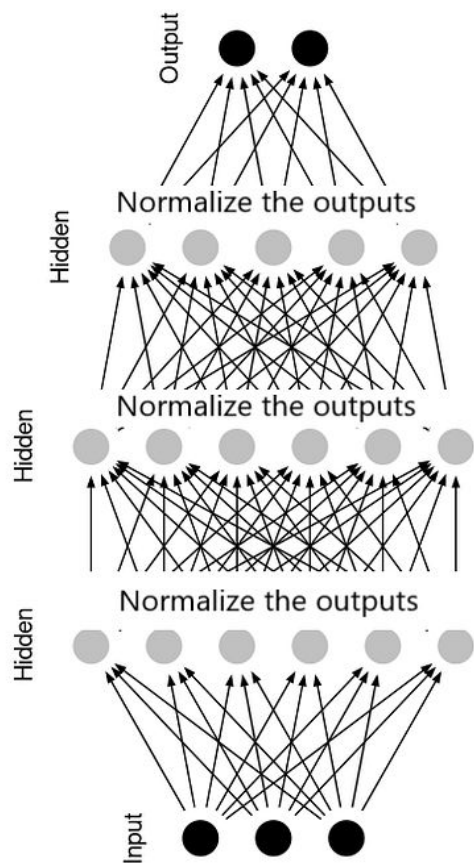
<https://www.analyticsvidhya.com/blog/2020/10/multivariate-multi-step-time-series-forecasting-using-stacked-lstm-sequence-to-sequence-autoencoder-in-tensorflow-2-0-keras/>

Vanishing/Exploding gradient → Residual Connections



A. Zhang+, https://d2l.ai/chapter_convolutional-modern/resnet.html

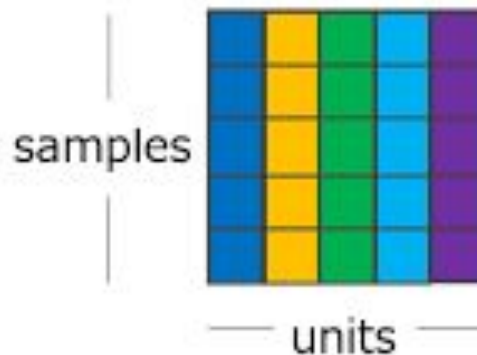
Covariate Shift → Batch / Layer Normalization



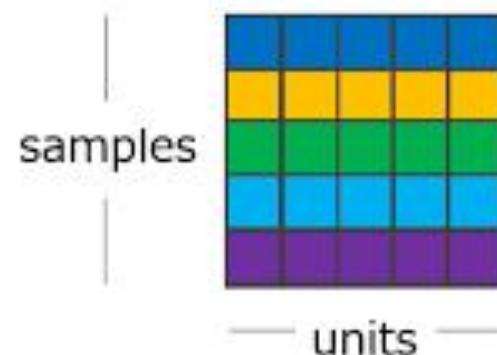
Normalize the outputs



Batch Normalization



Layer Normalization

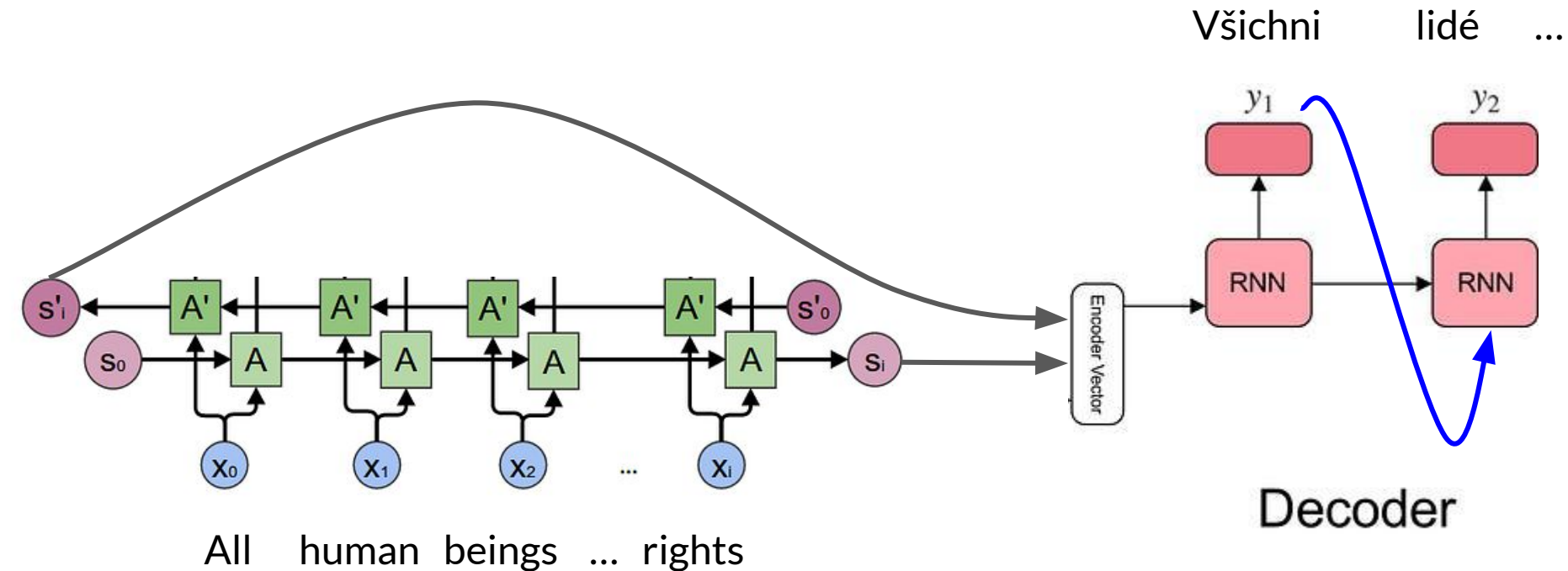


I. Rajagopal,

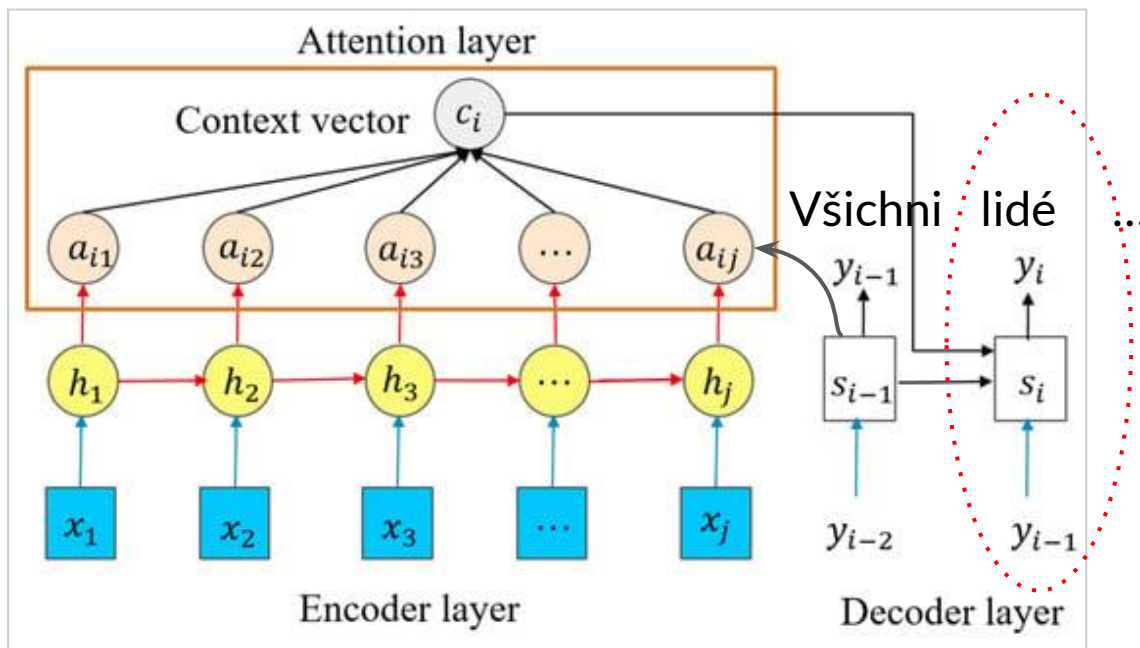
<https://medium.com/@ilango100/batch-normalization-speed-up-neural-network-training-245e39a62f85>

<https://ai-pool.com/a/s/normalization-in-deep-learning>

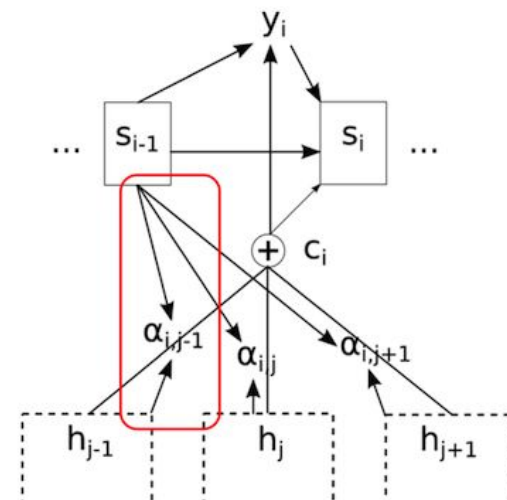
Encoding mostly the end → Bi-directional RNN



Selectively look at previous words' hidden states → Attention



All human beings ... rights

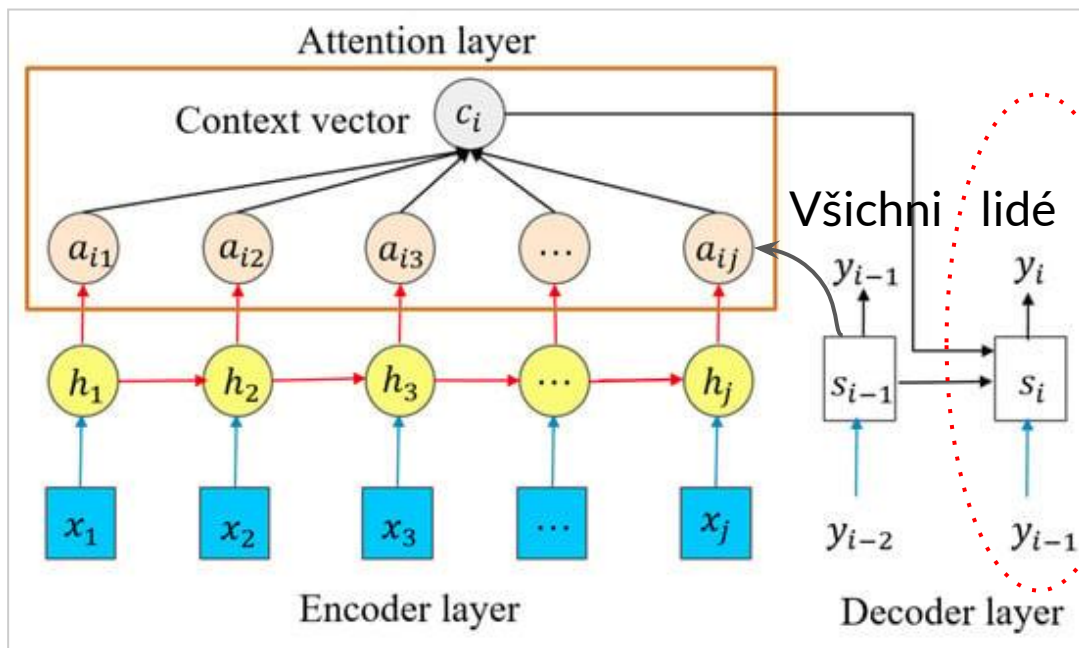


$$a_{i,j} = \text{softmax}(e_{i,j})$$
$$e_{i,j} = v^T \tanh(W_1 h_i + W_2 s_{i-1})$$

D. Bahdanau+, https://humanbrain.gitbook.io/notes/natural-language-processing/bahdanau_attention

M. Yang+, <https://doi.org/10.3390/electronics10141657>

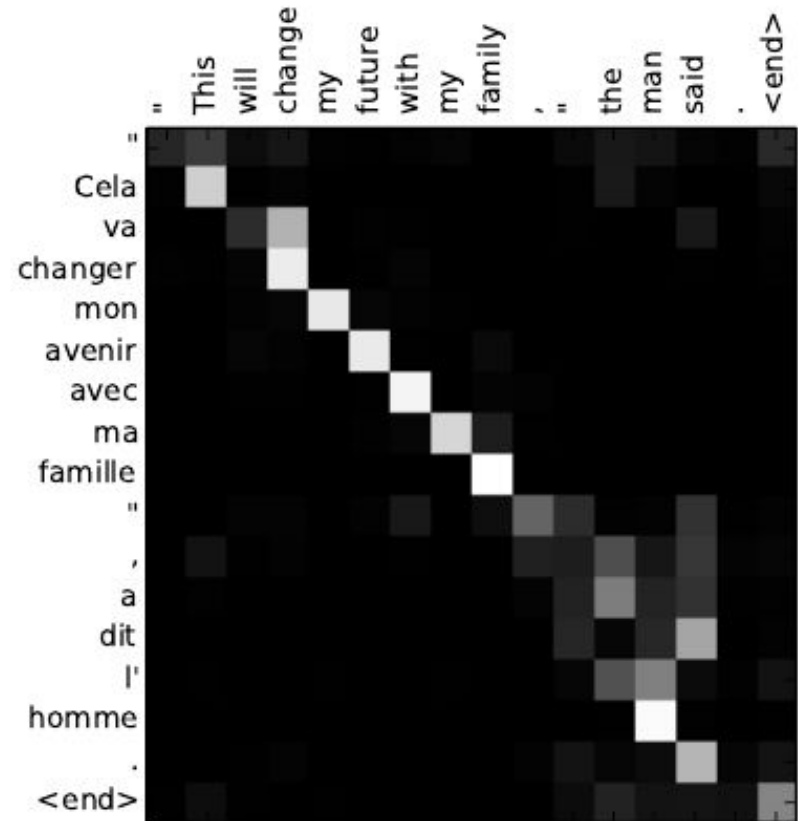
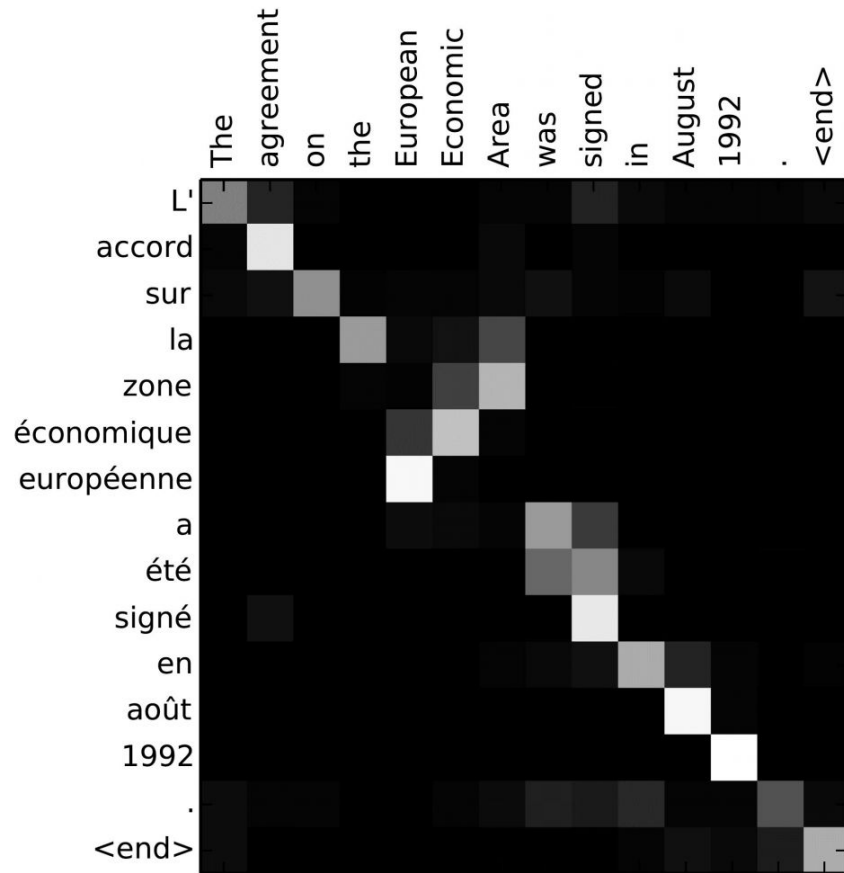
Selectively look at previous words' hidden states → Attention



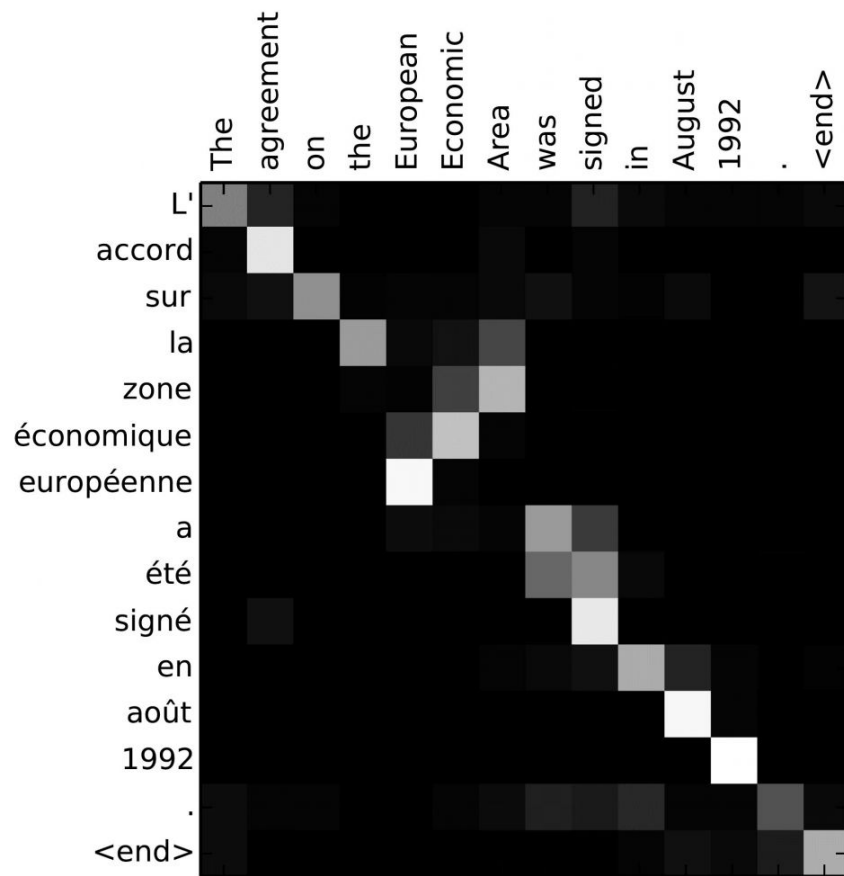
All human beings ... rights

- inspiration: word alignment
 - attention ~ soft alignment
- typically attends to word(s) being translated
 - “lidé” ~ “human beings”
- contextual word embs.
 - vs static word embeddings
 - hidden state ~ word representation in context
 - ELMo
- encoder final state useless
 - decode from a static state
 - use biRNN for hidden states

Visualisation of Attention



Visualisation of Attention



- this is single-head attention
 - typically attends to the word that is being translated
 - theoretically soft distribution
 - but typically very peaked
- also possible: multihead attention
 - one head attends to translated word
 - one head attends to preceding word
 - one head attends to sentence start
 - one head attends to sentence end
 - one head attends to main verb
 - one head attends to sentence subject
 - ...
 - learned (emergent), not enforced!

A range of useful components and concepts

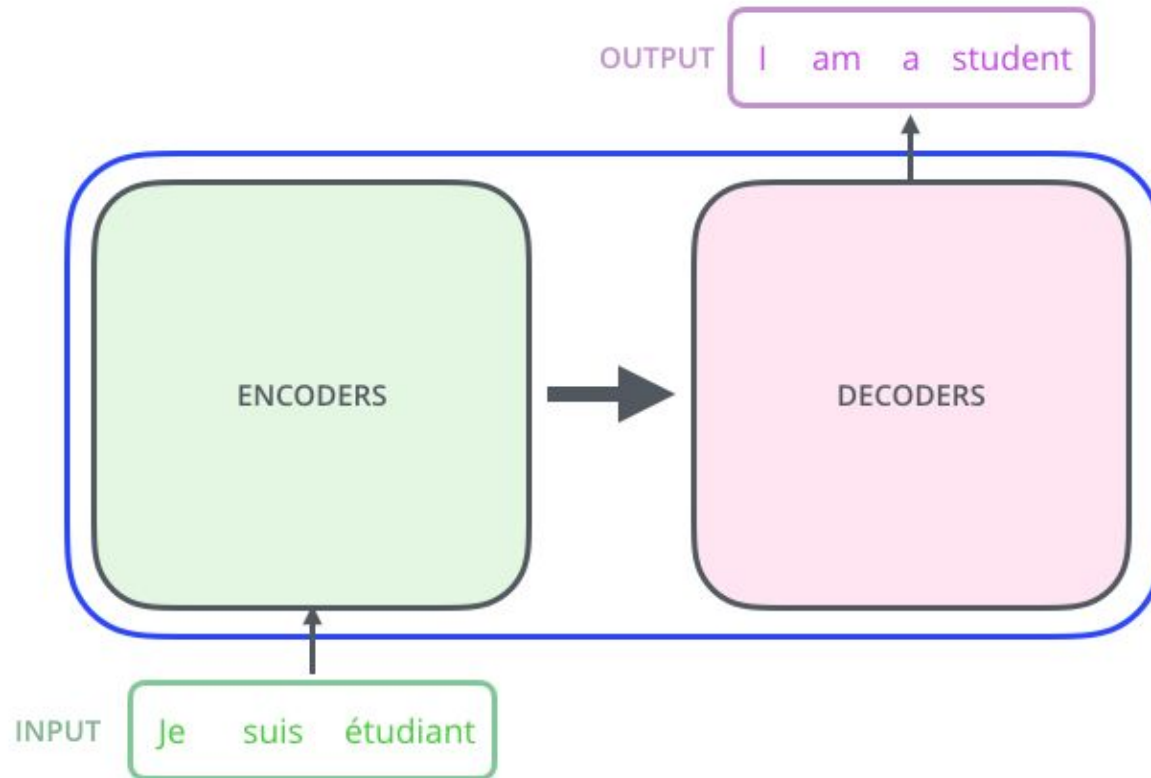
- NN cells (RNN, LSTM, GRU)
 - Note: CNNs possible but not common; sequence-to-label vs sequence-to-sequence
 - Note: Transformer is yet another alternative (to RNN or CNN) for sequence processing
- End-to-end
- Sequence-to-sequence
- Encoder-decoder
- Autoregressiveness
- Stacking
- Residual Connections (Highway Networks)
- Batch/Layer Normalization
- Directionality
- Attention
 - very powerful → key component of Transformer (“Attention is All You Need”)

The Transformer

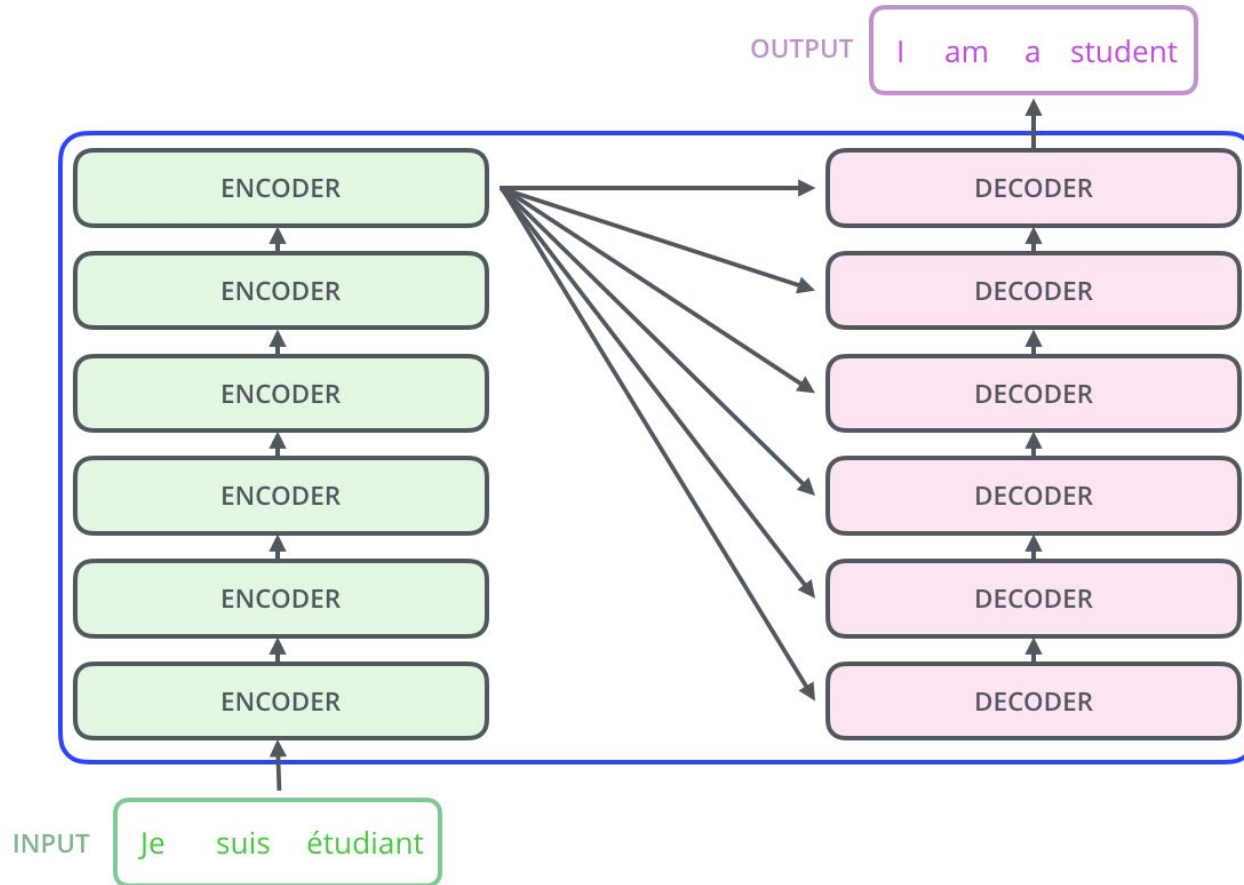
Transformer: end2end seq2seq



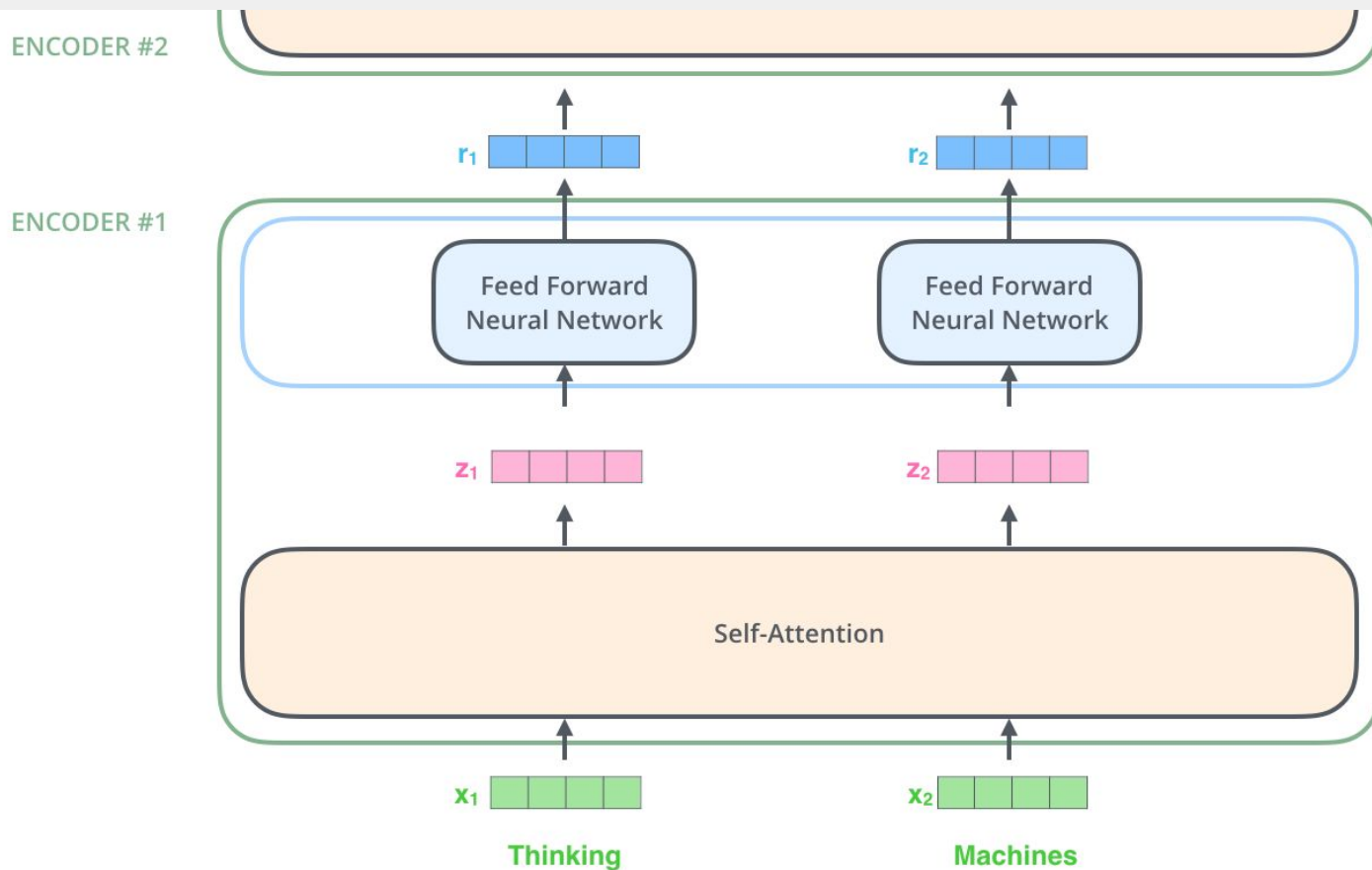
Transformer: encoder-decoder



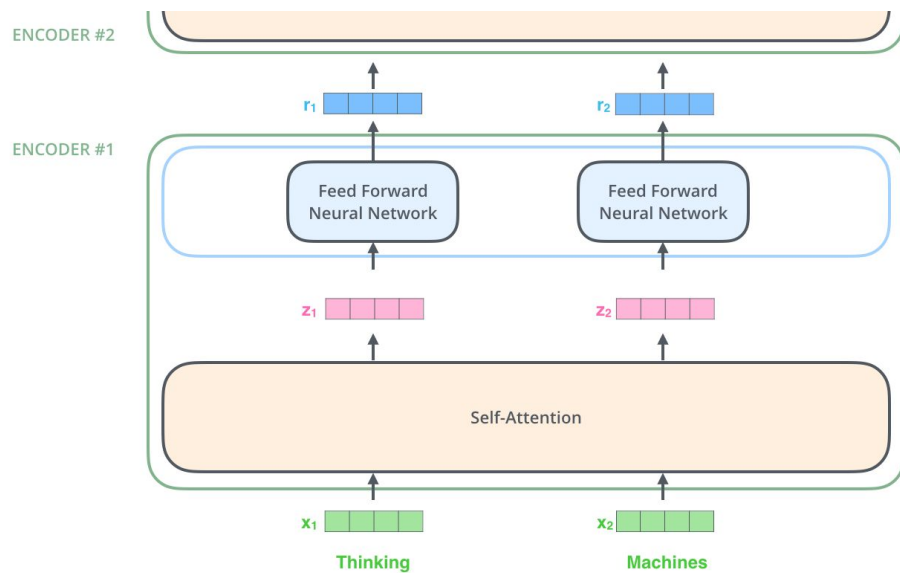
Transformer: stacking



Transformer: encoder; hidden state ~ contextual word repres.

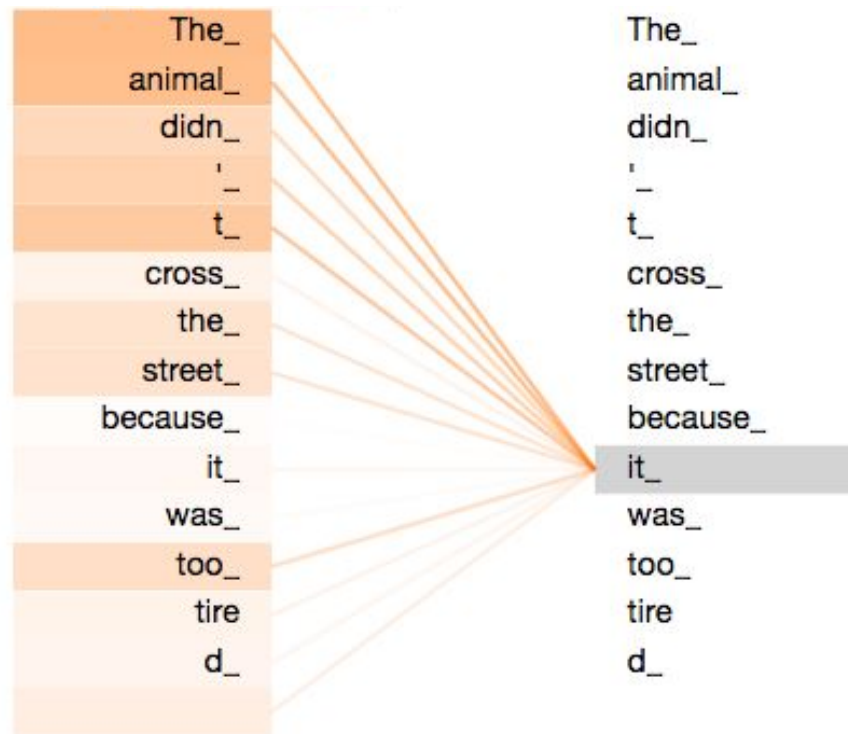


Transformer: self-attention (instead of RNN cell)

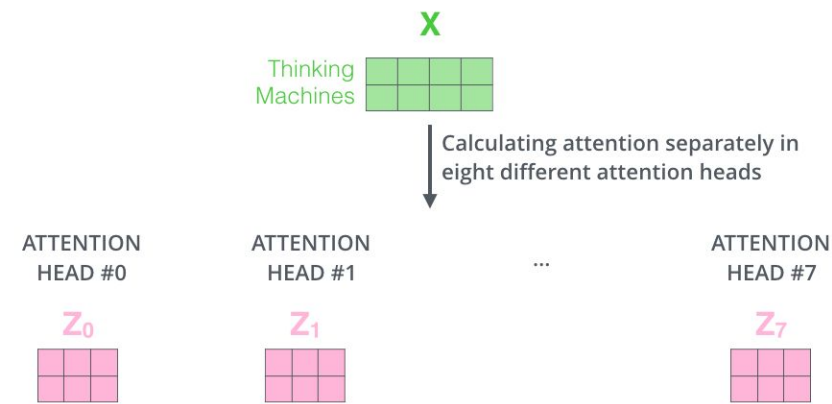


encoder input (x)

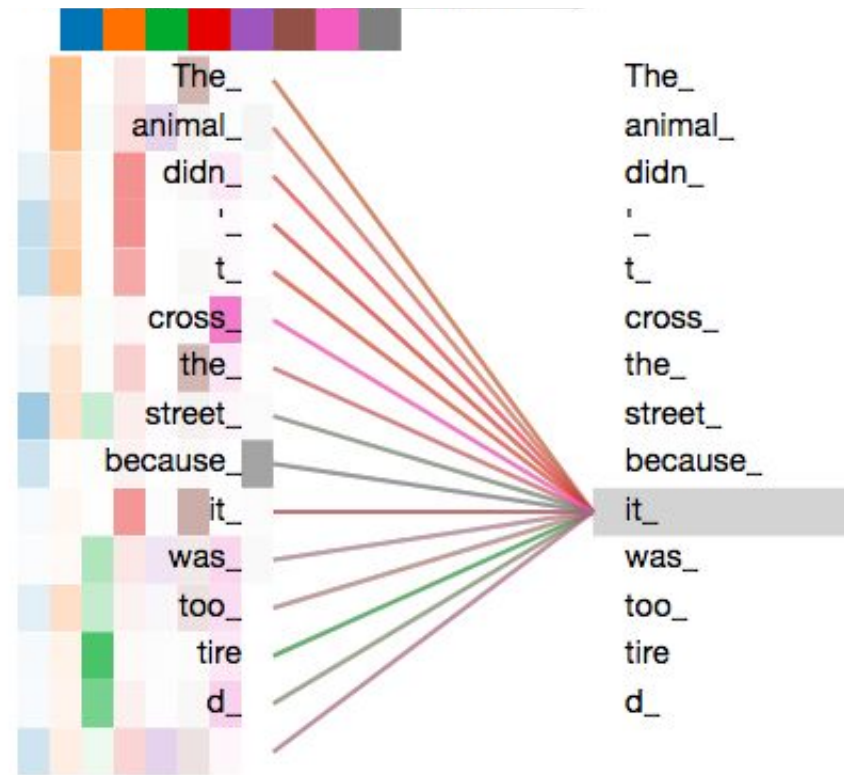
enc. internal state (z)



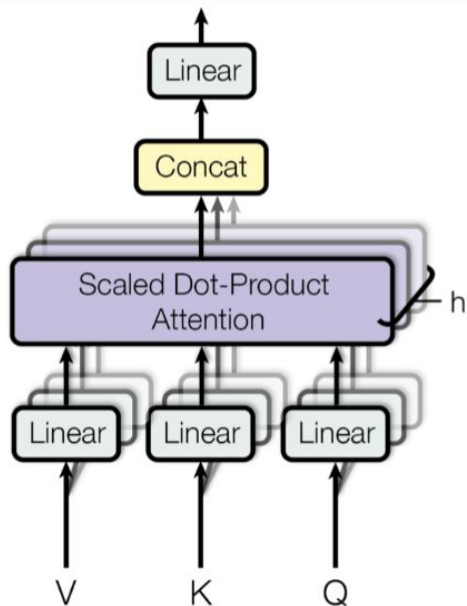
Transformer: multihead self-attention



1) Concatenate all the attention heads



Transformer: dot-product attention



Input

Thinking

Machines

Embedding

x_1

x_2

Queries

q_1

q_2



W^Q

Keys

k_1

k_2



W^K

Values

v_1

v_2

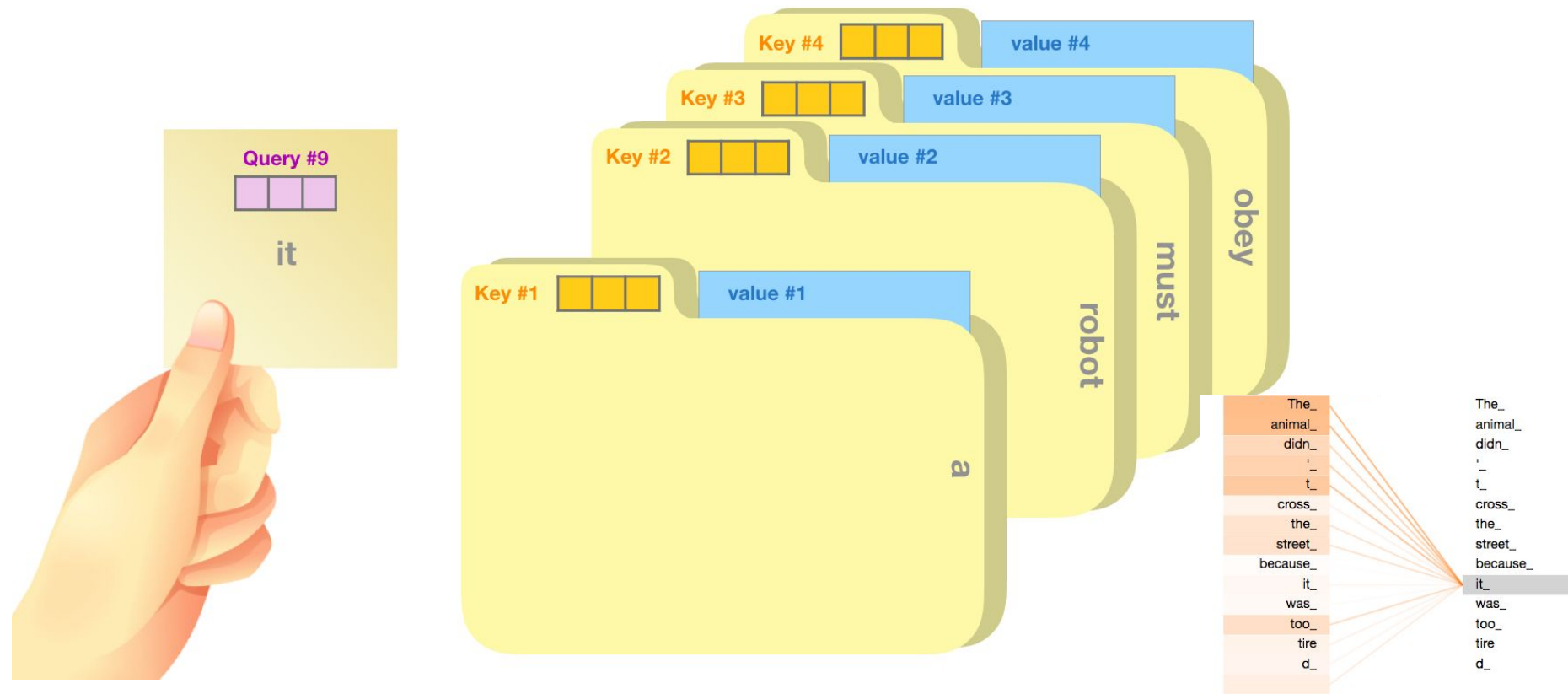


W^V

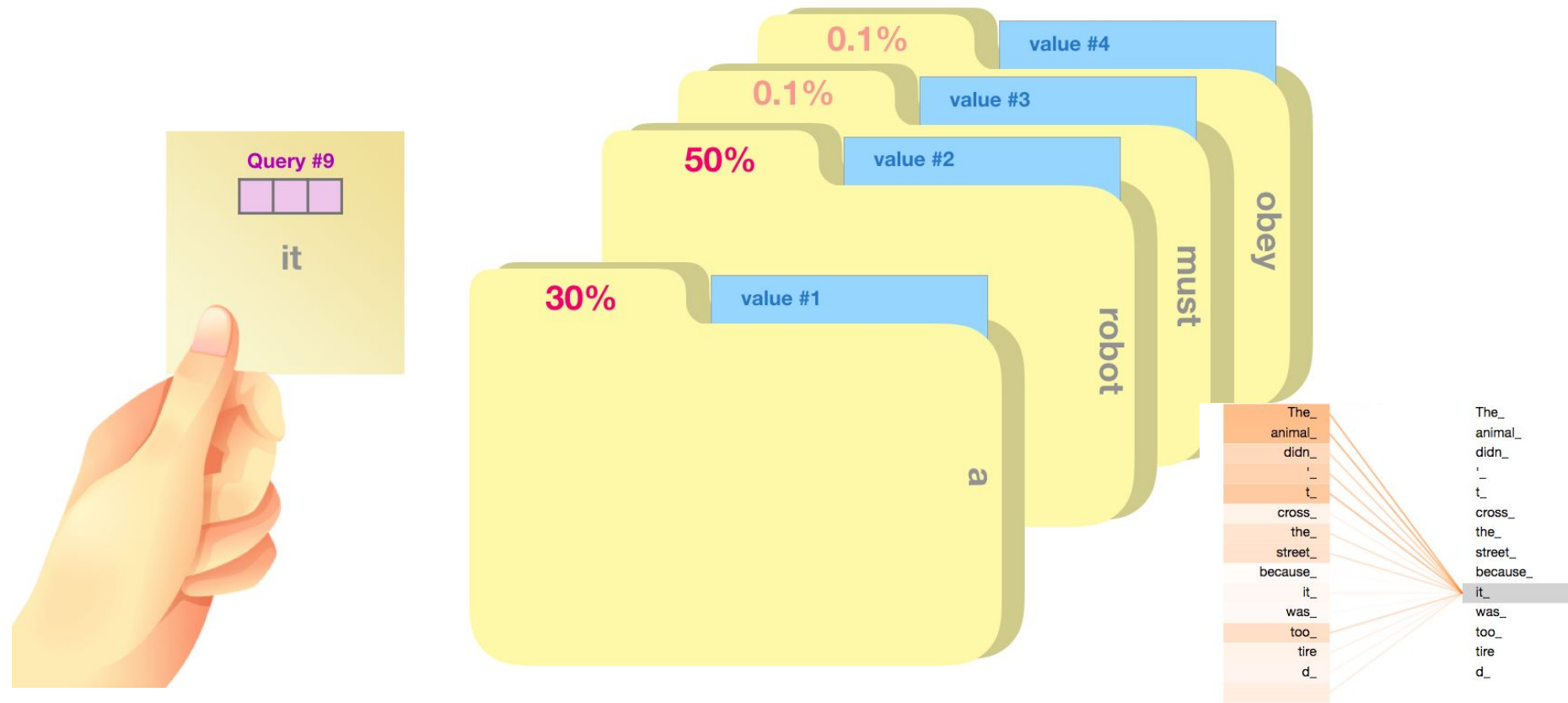
x_1

Thinking

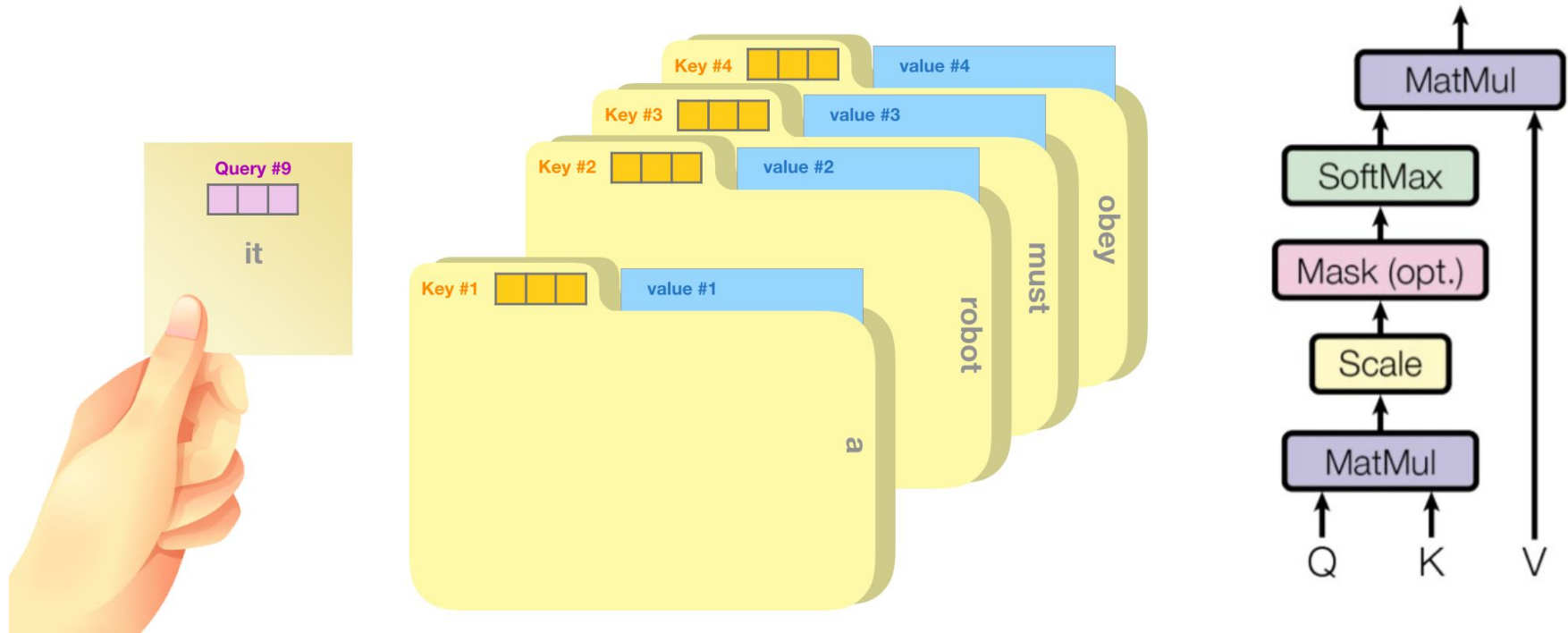
Transformer: dot-product attention



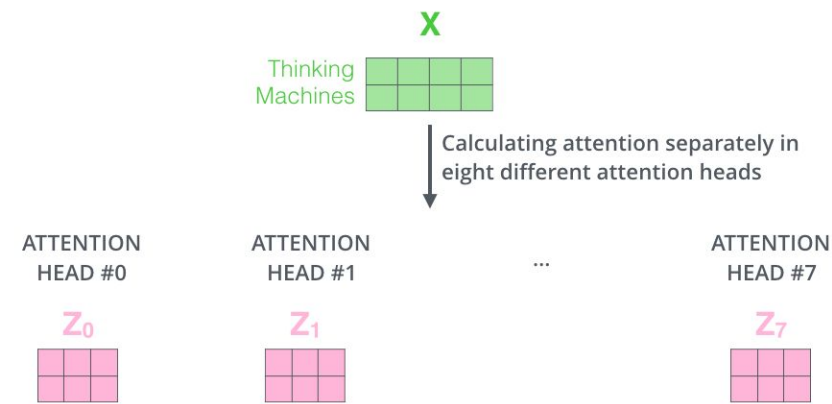
Transformer: dot-product attention



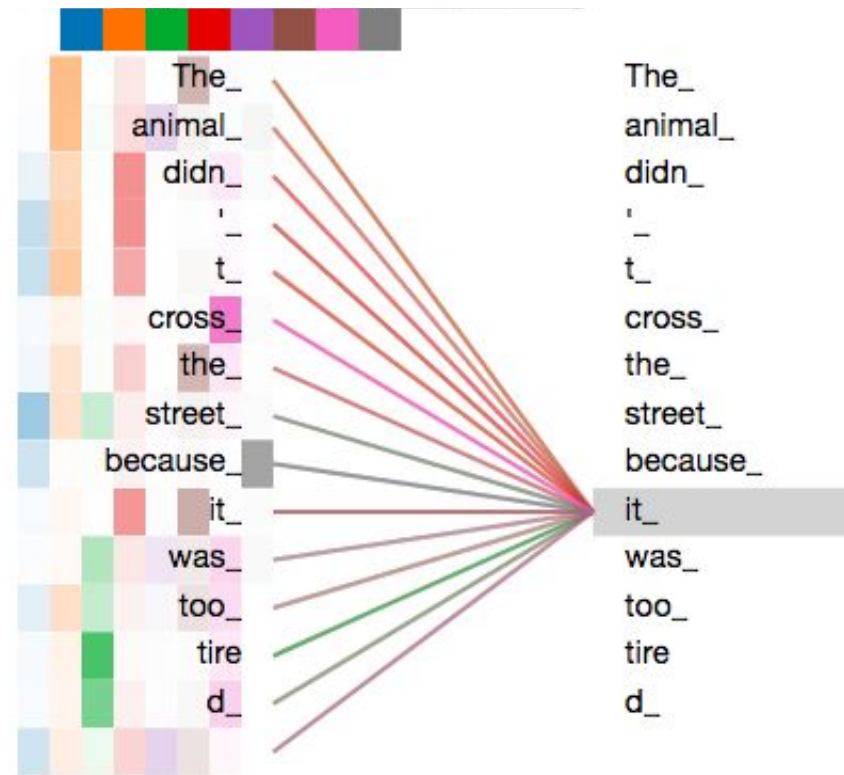
Transformer



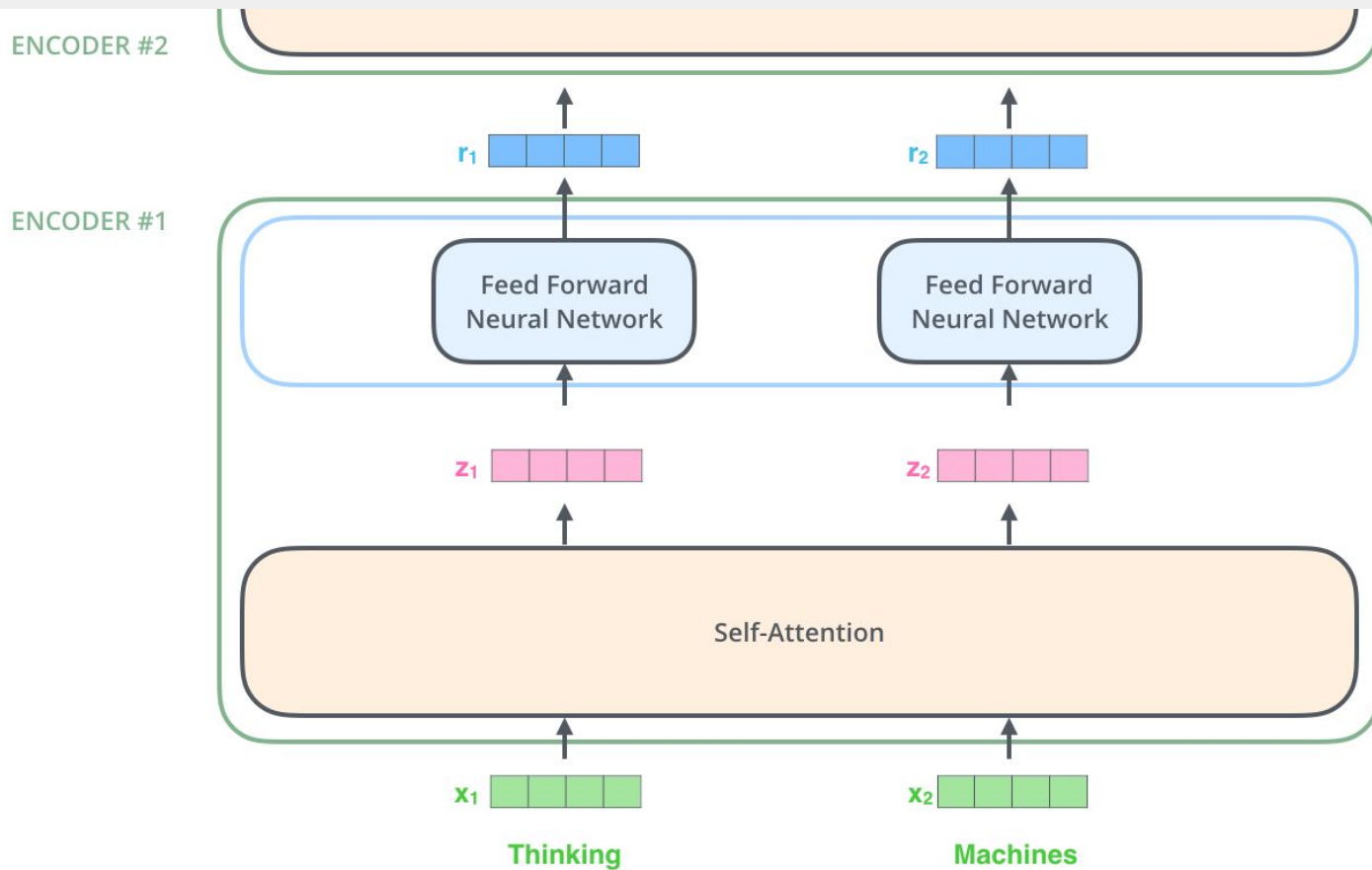
Transformer: multihead self-attention



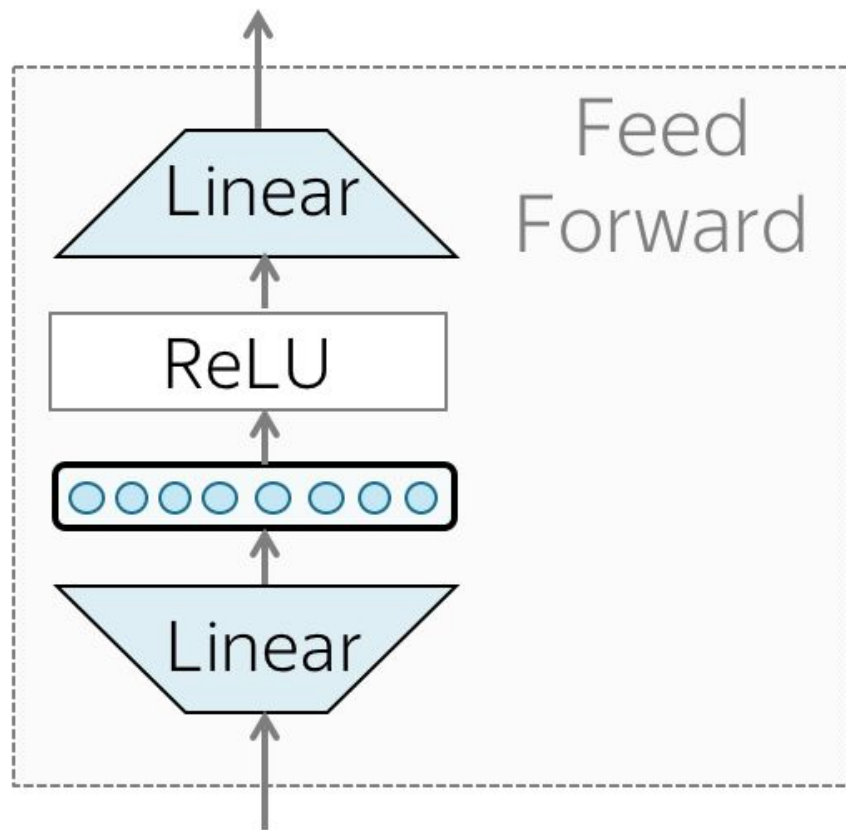
1) Concatenate all the attention heads



Transformer: encoder; hidden state ~ contextual word repres.



Transformer FFN: 1 hidden layer, “expand and contract” (4x)



Transformer FFN: 1 hidden layer, “expand and contract” (4x)

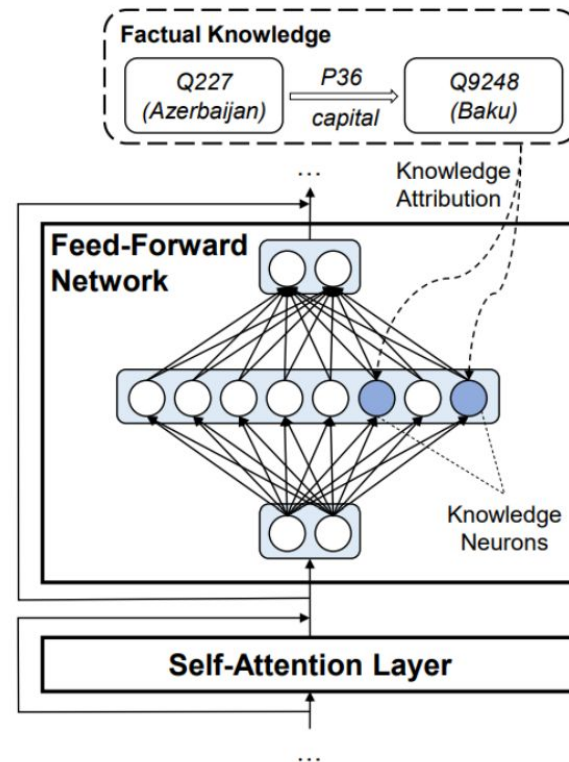
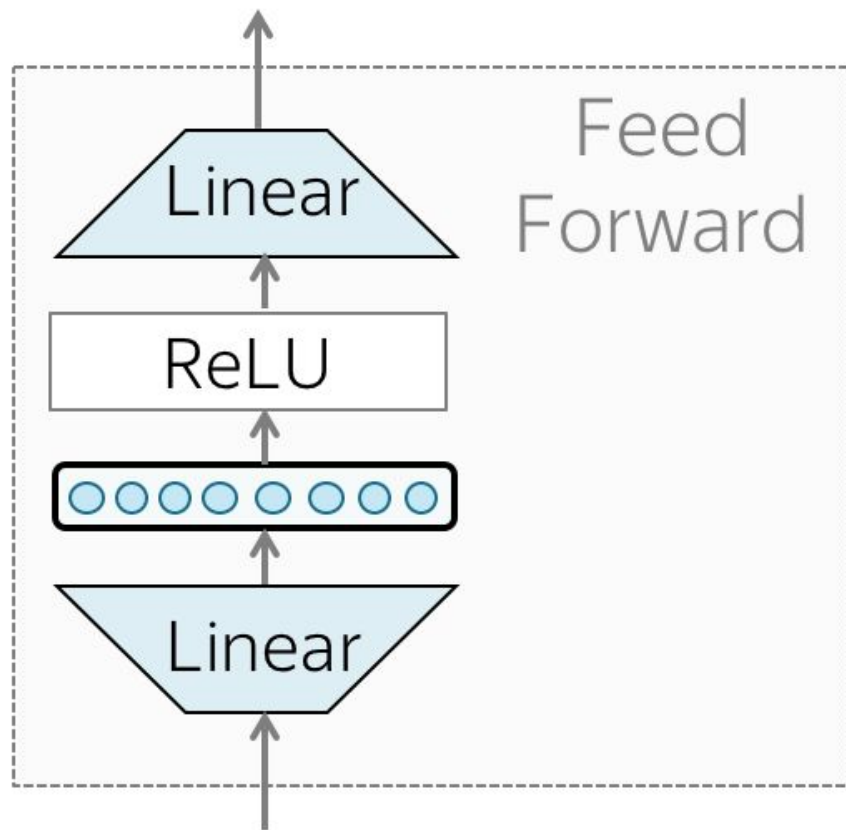
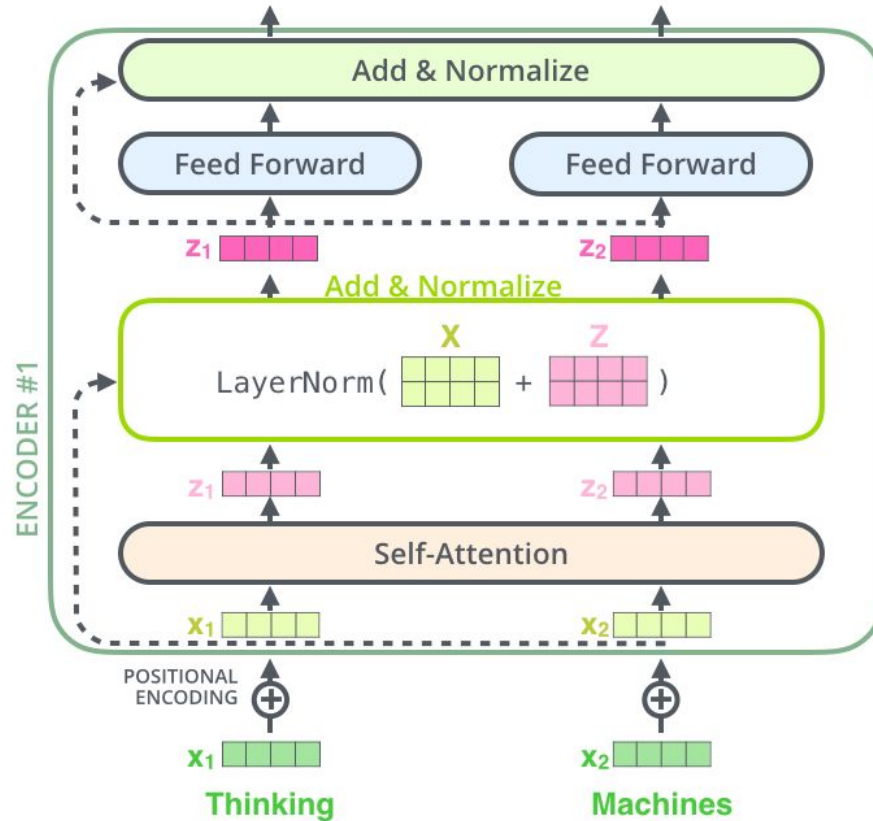
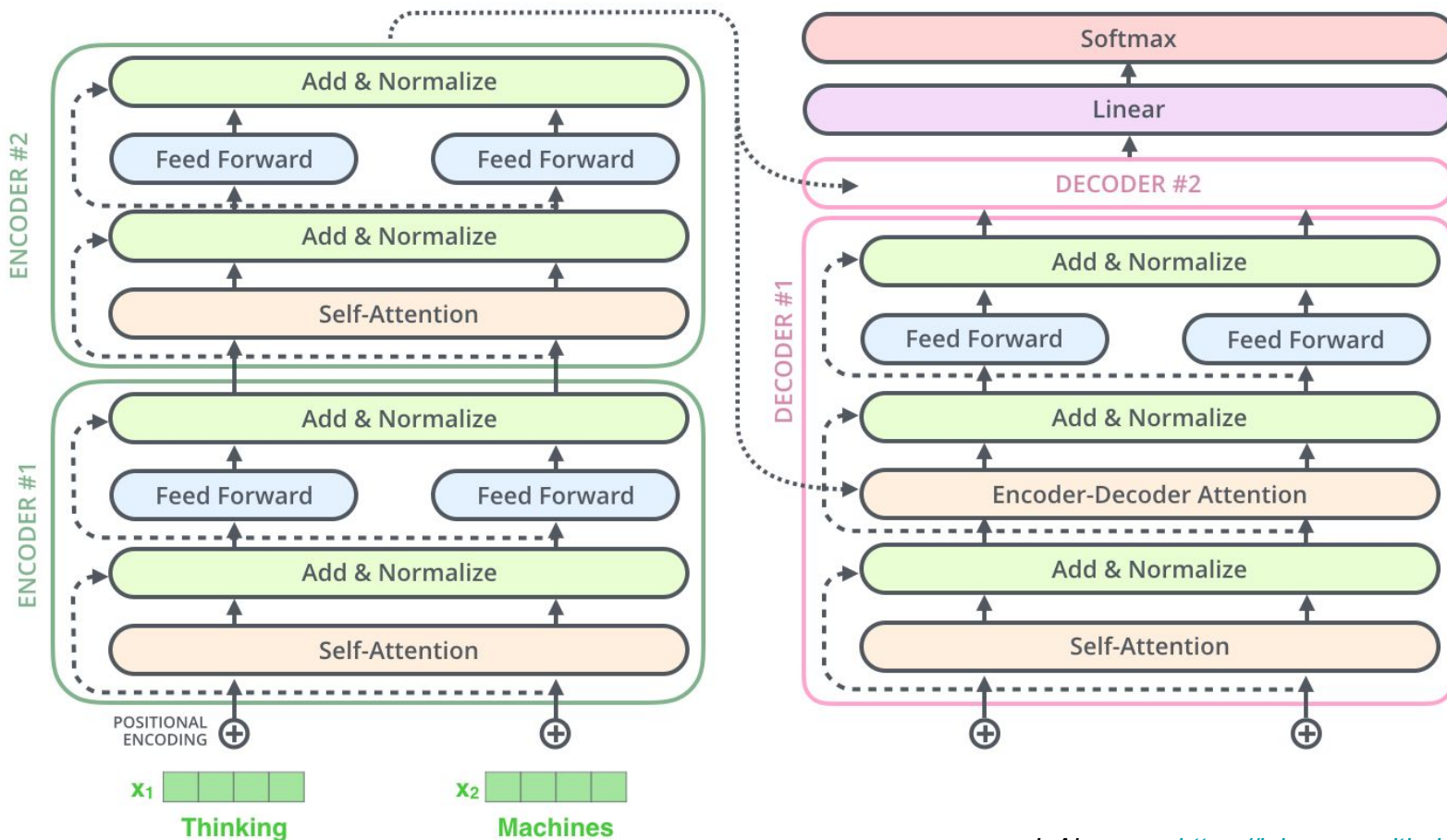


Figure 1: We aim to identify knowledge neurons correlated to a relational fact through knowledge attribution.

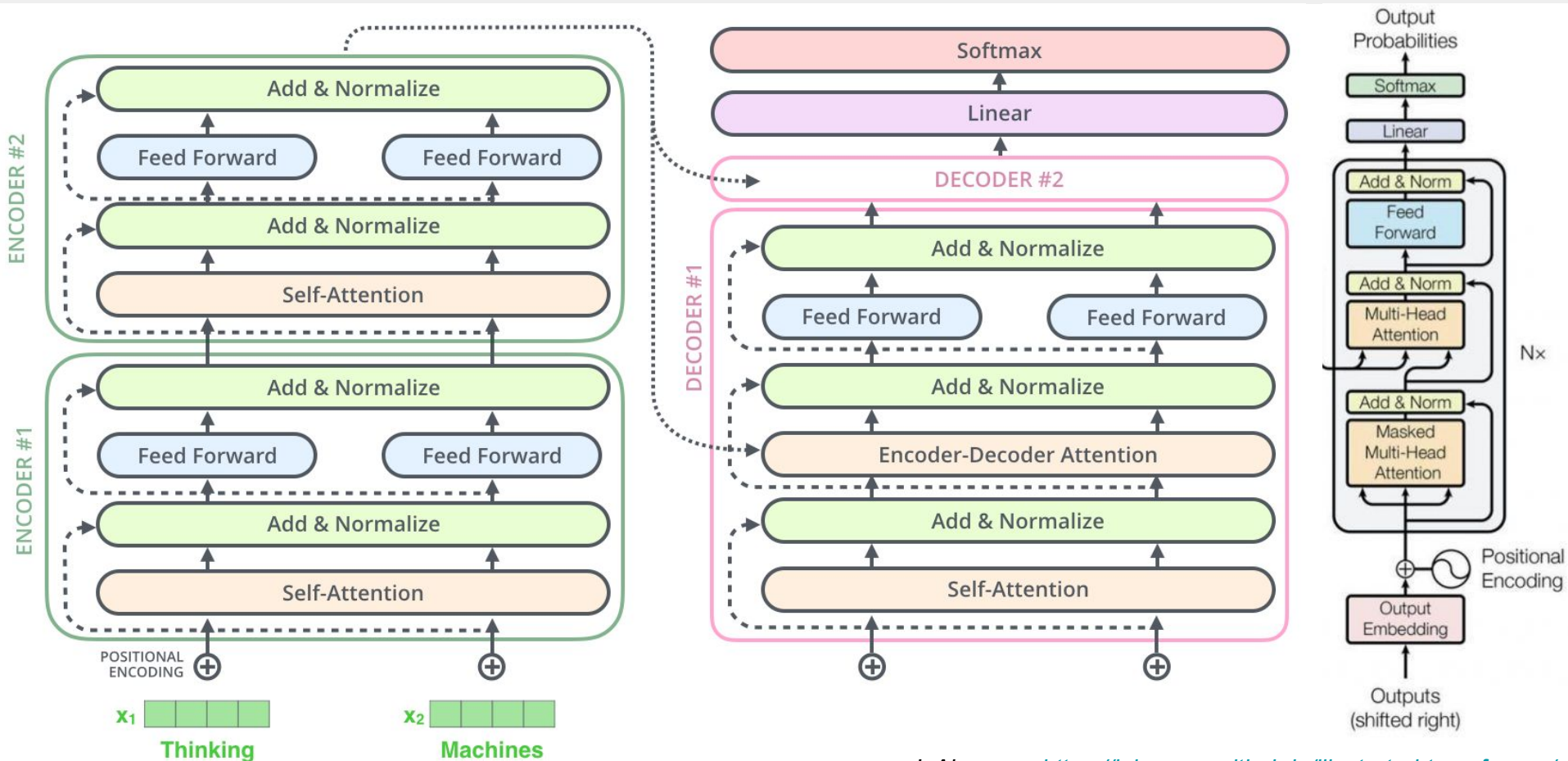
Transformer: full encoder block (residual connections, layer norm)



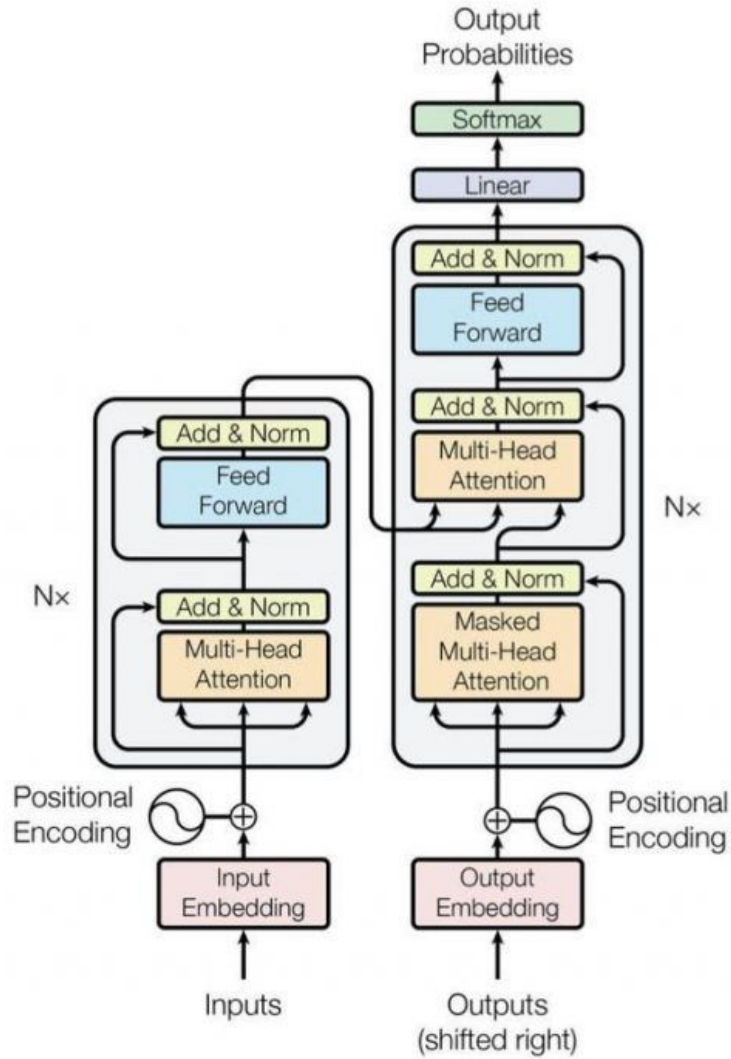
Transformer: encoder and decoder



Transformer: encoder and decoder



Transformer

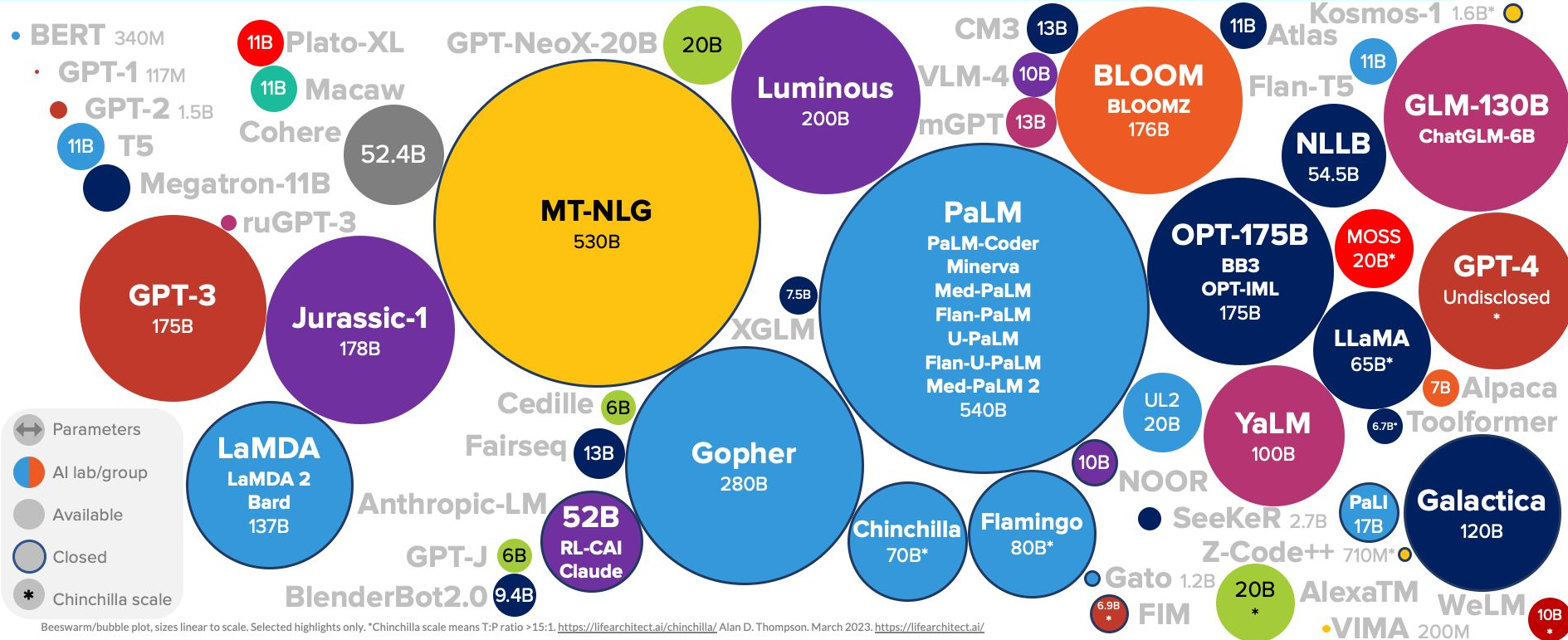


Transformers Models and Applications

Applications

- Sequence-to-sequence models
 - Machine Translation
 - parallel data (sentence → sentence' / paragraph → paragraph')
 - Text Summarisation
 - text → summary
 - ...
- Language Modelling
 - self-supervised (text → text)
 - next word prediction: *I woke up in the morning and went to the...*
 - autoregressive text generation
 - masked word prediction: *I woke up in the [MASK] and went to the bathroom.*
 - fill missing word, error correction
 - text representation (~ semantic tasks)
 - static word embeddings → contextual word embeddings → sentence embedding
 - text classification (part of speech, sentiment...)
 - ...

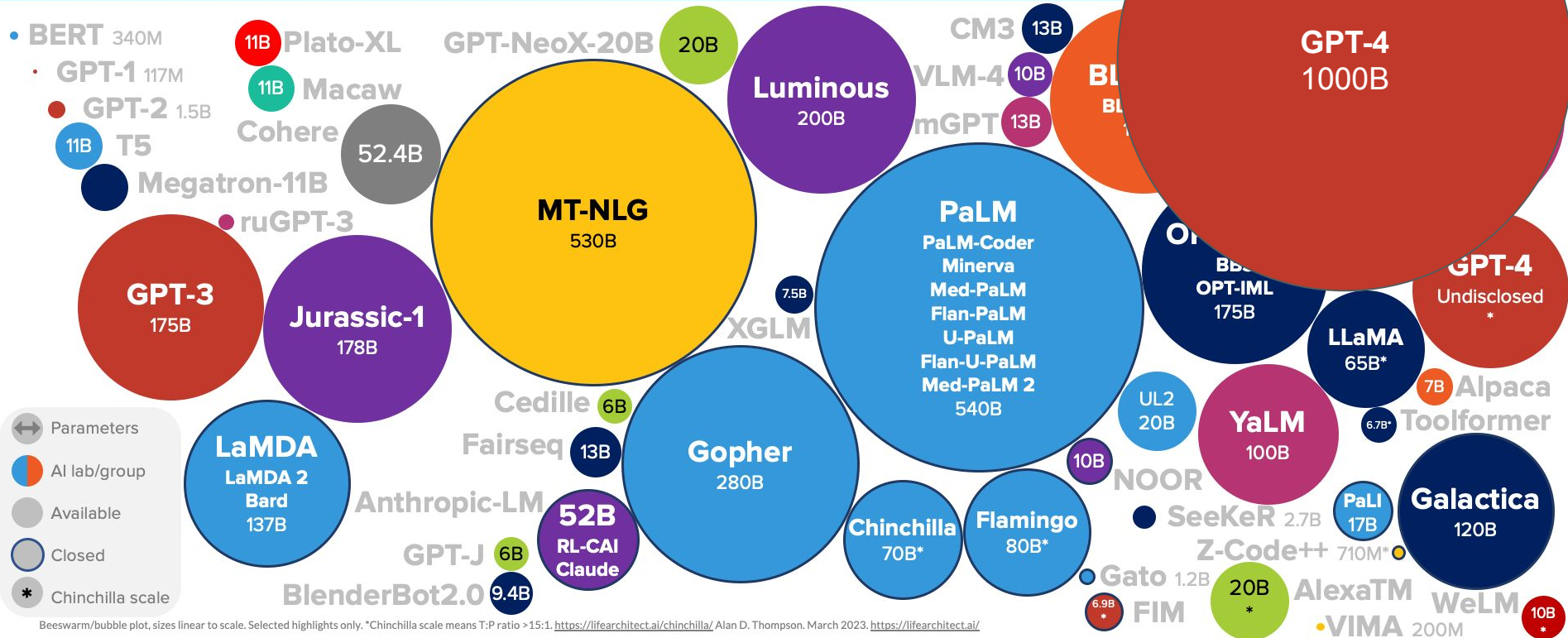
LANGUAGE MODEL SIZES TO MAR/2023



Beeswarm/bubble plot, sizes linear to scale. Selected highlights only. *Chinchilla scale means T:P ratio > 15:1. <https://lifearchitect.ai/chinchilla/> Alan D. Thompson, March 2023. <https://lifearchitect.ai/>



LANGUAGE MODEL SIZES TO MAR/2023



Beeswarm/bubble plot, sizes linear to scale. Selected highlights only. *Chinchilla scale means T:P ratio >15:1. <https://lifearchitect.ai/chinchilla/> Alan D. Thompson, March 2023. <https://lifearchitect.ai/>



Notable language models

- **masked**
 - BERT, Multilingual BERT, RoBERTa, XLM-RoBERTa, BART
- **generative**
 - OpenAI: GPT models (GPT-1 ... GPT-4, ChatGPT)
 - Meta: OPT, LLaMA → Alpaca, Vikuna, gpt4all
 - academia: HPLT, OpenGPT-X, OpenAssistant (all in progress)
- **model distillation (e.g. DistilBERT)**
 - teacher-student: train a large model, teach a small model to simulate the large model
- **going beyond naturally occurring text (i.e. crawling the internet) in GPT-3.5**
 - instructions (InstructGPT): “Summarize the following text: ...”
 - source codes (Copilot): “Write a Python script to sort an array...”
 - reinforcement learning with human feedback (RLHF): 👍👎 “correct answer is: ...”
- **going beyond text → multimodal (images), external tools, embodiment (robots)**
 - GPT-4, Retrieval-enriched LMs, Bing AI Bot, ToolFormer, ChatGPT plugins, PaLM-E

(Some) Text Generation Free Online Demos

- our demo
 - THEaiTRobot demo
 - GPT-2, fine-tuned for theatre script generation
 - <https://theaitre.com/demo>
- OpenAI
 - OpenAI Playground
 - based on GPT-3, various variants
 - <https://beta.openai.com/playground>
 - ChatGPT
 - released 30th Nov 2022
 - based on GPT-3.5 and [InstructGPT](#), trained with humans in the loop
 - <https://chat.openai.com>

CAN A ROBOT WRITE A THEATRE PLAY?

SAMPLE SCRIPT

Scene Polonius speaks to the king. Enter Hamlet.

Polonius I hear him coming; let's hide, sir.

Hamlet To be or not to be; that is the question.



GENERATE FURTHER

INPUT YOUR OWN SCRIPT

AI will generate a continuation

Scene Write here a description of the starting situation.

Name of first character Write here what the first character says.

Name of second character Write here what the second character says

GENERATE FURTHER

Scene

Helen's salon. In the room on the left, Helen plays the piano. Dominus paces the room, Dr. Gall looks out the window, and Alquist sits off to one side in a lounge chair with his face covered by his hands.

Dr. Gall

Heavens, there's more!

Domin

Robots?

Dr. Gall

That's right! Now what?

Alquist

I was afraid it would get worse.

Dominus

I am sorry, Dr. Gall. I'm afraid I have no choice but to kill you.



theaitre.com/demo



The screenshot shows the OpenAI Playground interface in a browser. The address bar shows the URL `beta.openai.com/playground`. The page title is "Playground - OpenAI API". The interface includes a search icon, a help icon, and a user profile labeled "Personal". Below the navigation bar, there is a "Playground" section with a "Load a preset..." dropdown menu and buttons for "Save", "View code", "Share", and a menu icon. The main text area contains the following text:

Martin Platek, František Mráz, Dana Pardubská and Daniel Prusa: On Pumping RP-automata controlled by complete LRG(\cent, \$)-grammars

Abstract:

At the bottom of the text area, there is a "Submit" button with a tooltip that says "Submit Ctrl Enter". Below the text area, there is a row of controls: a "Submit" button, a refresh icon, a redo icon, a undo icon, and a token count of "47".

The screenshot shows the OpenAI Playground interface. At the top, the browser address bar displays "beta.openai.com/playground". The page title is "Playground - OpenAI API". The OpenAI logo is on the left, and "Help" and "Personal" are on the right. Below the header, there's a "Playground" section with a "Load a preset..." dropdown, "Save", "View code", "Share", and a menu icon. The main content area contains a text snippet with a green highlight. At the bottom, there's a "Regenerate" button with "Ctrl ↑ Enter" instructions, and a row of action buttons: "Submit", a refresh icon, a redo icon, a undo icon, a quote icon, a thumbs up icon, and a "47" badge.

Playground - OpenAI API x

beta.openai.com/playground

Help Personal

Playground Load a preset... Save View code Share ...

Martin [Platek](#), [František Mráz](#), Dana [Pardubská](#) and Daniel [Prusa](#): On Pumping RP-automata controlled by complete [LRG\(\$\cent, \S\$ \)-grammars](#)

Abstract:

We prove that one-way [RM\(\$\cent, \S\$ \)-automata](#) over infinite alphabets can be defined by complete [LRG\(\$\cent, \S\$ \)-grammars](#) with a one-way pumping lemma. We also show that there is a close connection between pumping of one-way [RM\(\$\cent, \S\$ \)-automata](#) and pumping of one-way deterministic context-free languages, and we prove that one-way deterministic context-free languages can be defined by complete [LRG\(\$\cent, \S\$ \)-grammars](#) with a one-way pumping lemma.

To appear in the Proceeding of the 12th International Conference on Implementation and Application of Automata, [Springer-Verlag](#), Lecture Notes in Computer Science ([LNCS](#)) 3181, 2004.

Regenerate Ctrl ↑ Enter

Submit ↺ ↻ ↶ ↷ 👍 47



Please write the abstract for the paper "Analyzing Item Popularity Bias of Music Recommender Systems: Are Different Genders Equally Affected?"



I'm sorry, I couldn't find a paper with that exact title. Recommender systems research is a broad field, and there could be papers with similar titles.



It is also possible that the paper you are looking for hasn't been published yet, or that it is not publicly available.

However, a possible abstract for a paper with that title could be:

"Music recommender systems have become an integral part of many music streaming platforms, yet they have been shown to suffer from various biases, such as item popularity bias. This bias leads to the recommendation of popular items more frequently than less popular items. In this paper, we analyze the extent to which item popularity bias affects the recommendations made by music recommender systems for different genders. Using a dataset of music listening history, we compare the popularity of items recommended to male and female users and investigate whether there are significant differences in the level of bias. Our results show that both male and female users are affected by item popularity bias, but to different degrees. These findings highlight the importance of considering and addressing bias in music recommender systems, particularly in terms of gender."

This is a hypothetical abstract for a paper that might be out there, and I apologize if this is not the one you were looking for.



I'm sorry, I couldn't find a paper with that exact title. Recommender systems research is a broad field, and there could be papers with similar titles.



It is also possible that the paper you are looking for hasn't been published yet, or that it is not publicly available.

However, a possible abstract for a paper with that title could be:

"Music recommender systems have become an integral part of many music streaming platforms, yet they have been shown to suffer from various biases, such as item popularity bias. This bias leads to the recommendation of popular items more frequently than less popular items. In this paper, we analyze the extent to which item popularity bias affects the recommendations made by music recommender systems for different genders. Using a dataset of music listening history, we compare the popularity of items recommended to male and female users and investigate whether there are significant differences in the level of bias. Our results show that both male and female users are affected by item popularity bias, but to different degrees. These findings highlight the importance of considering and addressing bias in music recommender systems, particularly in terms of gender."

This is a hypothetical abstract for a paper that might be out there, and I apologize if this is not the one you were looking for.

Some language model limitations

- **Based on available data**
 - good for frequent text types, bad for rare types/topics/events/styles/genres...
 - biases, stereotypes, racism, sexism, clichés...
 - data hungry, needs very large training data (e.g. compared to people)
- **Lacks knowledge of knowledge**
 - intermixed language knowledge (probable text) and factual knowledge (true facts)
 - limited introspection/confidence estimation
 - research on interpretability and explainability
- **Limited channels**
 - text, limited multimodality, limited real-world experience, lacks human experience
- **Technical limitations**
 - Limited context ... GPT-4: window of up to 32k tokens (dozens of pages)
 - Bad notion of time (sees texts from all time points at once, no metadata)
 - Bad representation of numbers and maths
- **Understanding language? (General) Artificial Intelligence?**
 - ?

GPT et al.: Generating Texts with Transformer-Based LLMs

- Machine Translation, Language Modelling
- RNN, Attention, Transformer
- Large Language Models

These slides available at bit.ly/theaitre-sui

<https://www.theaitre.com/>